



Europäisches
Patentamt

European
Patent Office

Office européen
des brevets

PCT/EP 03/501903
03/006243

REC'D 08 APR 2003

WIPO

PCT

Bescheinigung

Certificate

Attestation

Die angehefteten Unterla-
gen stimmen mit der
ursprünglich eingereichten
Fassung der auf dem näch-
sten Blatt bezeichneten
europäischen Patentanmel-
dung überein.

The attached documents
are exact copies of the
European patent application
described on the following
page, as originally filed.

Les documents fixés à
cette attestation sont
conformes à la version
initialement déposée de
la demande de brevet
européen spécifiée à la
page suivante.

Patentanmeldung Nr. Patent application No. Demande de brevet n°

02001331.4

PRIORITY DOCUMENT

SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)

Der Präsident des Europäischen Patentamts;
Im Auftrag

For the President of the European Patent Office

Le Président de l'Office européen des brevets
p.o.

R C van Dijk

DEN HAAG, DEN
THE HAGUE,
LA HAYE, LE

31/03/03



Europäisches
Patentamt

European
Patent Office

Office européen
des brevets

**Blatt 2 der Bescheinigung
Sheet 2 of the certificate
Page 2 de l'attestation**

Anmeldung Nr.:
Application no.: 02001331.4
Demande n°:

Anmeldetag:
Date of filing: 18/01/02
Date de dépôt:

Anmelder:
Applicant(s):
Demandeur(s):
Pact Informationstechnologie GmbH
80939 Munchen
GERMANY

Bezeichnung der Erfindung:
Title of the invention:
Titre de l'invention:
Method and device

In Anspruch genommene Priorität(en) / Priority(ies) claimed / Priorité(s) revendiquée(s)

Staat:
State:
Pays:

Tag:
Date:
Date:

Aktenzeichen:
File no.
Numéro de dépôt:

Internationale Patentklassifikation:
International Patent classification:
Classification internationale des brevets:
/

Am Anmeldetag benannte Vertragsstaaten:
Contracting states designated at date of filing:
Etats contractants désignés lors du dépôt:

AT/BE/CH/CY/DE/DK/ES/FI/FR/GB/GR/IE/IT/LI/LU/MC/NL/PT/SE/TR

Bemerkungen:
Remarks:
Remarques:

THE APPLICANT'S ADDRESS AT THE TIME OF FILING OF THE APPLICATION WAS:
Leopoldstra e 236 -80807 Munich, Germany
THE REGISTRATION OF THE CHANGE HAS TAKEN EFFECT ON 19.04.02

Akte: PACT27/EP

Applicant: Pact Informationstechnologie GmbH
Leopoldstr. 236
80807 München

REpresentative:

European Patent Attorney
Claus Peter Pietruk
Heinrich-Lilienfein-Weg 5
D-76229 Karlsruhe

European Patent Application

The present invention relates to the subject matter claimed and hence refers to a method and a device for compiling programs for a reconfigurable device.

Reconfigurable devices are well-known. They include systolic arrays, neuronal networks, Multiprocessor systems, Prozessoren comprising a plurality of ALU and/or logic cells, crossbar-switches, as well as FPGAs, DPGAs, XPUTERS, asf.. Reference is being made to DE 44 16 881 A1, DE 197 81 412 A1, DE 197 81 483 A1, DE 196 54 846 A1, DE 196 54 593 A1, DE 197 04 044.6 A1, DE 198 80 129 A1, DE 198 61 088 A1, DE 199 80 312 A1, PCT/DE 00/01869, DE 100 36 627 A1, DE 100 28 397 A1, DE 101 10 530 A1, DE 101 11 014 A1, PCT/EP 00/10516, EP 01 102 674 A1, DE 198 80 128 A1, DE 101 39 170 A1, DE 198 09 640 A1, DE 199 26 538.0 A1, DE 100 50 442 A1 the full disclosure of which is incorporated herein for purposes of reference.

Furthermore, reference is being made to devices and methods as known from US PS 6,311,200; US PS 6,298,472; US PS 6,288,566; US PS 6,282,627; US PS 6,243,808 issued to Chameleon systems INC, USA noting that the disclosure of the present application

Akte: PACT27/EP

is pertinent in at least some aspects to some of the devices disclosed therein.

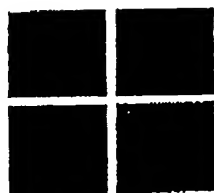
The invention will now be described by the following papers which are part of the present application.

Akte: PACT27/EP

Priority:

DE 101 39 170.6
 DE 101 42 903.7
 DE 101 44 732.9
 DE 101 45 792.8
 DE 10 154 260.7
 PCT/EP 01/088534
 PCT/EP 01/06703
 DE 101 10 530.4
 DE 101 11 014.5
 DE 101 46 132.1
 DE 101 29 237.6
 EP 01 115 021.6
 DE 101 42 494.4
 DE 101 45 795.2
 DE 101 42 904.5
 DE 101 44 733.7
 DE 101 54 259.3
 PCT/EP 01/11299
 PCT/EP 01/11593
 EP 01 102 674.7
 DE 101 35 210.7
 DE 101 35 211.5
 DE 101 42 231.8
 EP 01 129 923.7

Public



PACT

XPP Technologies

Vectorizing C Compiler XPP-VC

Introduction and Manual

© 2002 PACT XPP Technologies

Version 1.1

2002-1-11

1 Introduction

This document describes the PACT Vectorizing C Compiler XPP-VC which maps a C subset extended by port access functions to PACT's Native Mapping Language NML. A future extension of this compiler for a host-XPP hybrid system is described in Section 7.3.

XPP-VC uses the public domain SUIF compiler system. For installation instructions on both SUIF and XPP-VC, refer to the separately available installation notes.

2 General Approach

The XPP-VC implementation is based on the public domain SUIF compiler framework (cf. <http://suif.stanford.edu>). SUIF was chosen because it is easily extensible.

SUIF was extended with two passes: *partition* and *nmlgen*. The first pass, *partition*, tests if the program complies with the restrictions of the compiler (cf. Section 3.1) and performs a dependence analysis. It determines if a FOR-loop can be vectorized and annotates the syntax tree accordingly. In XPP-VC, vectorization means that loop iterations are overlapped and executed in a pipelined, parallel fashion. This technique is based on the *Pipeline Vectorization* method developed for reconfigurable architectures.¹ *partition* also completely unrolls inner program FOR-loops which are annotated by the user. All innermost loops (after unrolling) which can be vectorized are selected and annotated for pipeline synthesis.

nmlgen generates a control/dataflow graph for the program as follows. First, program data is allocated on the XPP Core. By default, *nmlgen* maps each program array to internal RAM blocks while scalar variables are stored in registers within the PAEs. If instructed by a pragma directive (cf. Section 3.2.2), arrays are mapped to external RAM. If it is large enough, an external RAM can hold several arrays.

Next, one ALU is allocated for each operator in the program (after loop unrolling, if applicable). The ALUs are connected according to the data-flow of the program. This data-driven execution of the operators automatically yields some instruction-level parallelism within a basic block of the program, but the basic blocks are normally executed in their original, sequential order, controlled by event signals. However, for generating more efficient XPP Core configurations, *nmlgen* generates pipelined operator networks for inner program loops which have been annotated for vectorization by *partition*. In other words, subsequent loop iterations are started before previous iterations have finished. Data packets flow continuously through the operator pipelines. By applying pipeline balancing techniques, maximum throughput is achieved. For many programs, additional performance gains are achieved by the complete loop unrolling transformation. Though unrolled loops require more XPP resources because individual PAEs are allocated for each loop iteration, they yield more parallelism and better exploitation of the XPP Core.

Finally, *nmlgen* outputs a self-contained NML file containing a module which implements the program on an XPP Core. The XPP IP parameters for the generated NML file are read from a configuration file, cf. Section 4. Thus the parameters can be easily changed. Obviously, large programs may produce NML files which cannot be placed and routed on a given XPP Core. Later XPP-VC releases will perform a temporal partitioning of C programs in order to overcome this limitation, cf. Section 7.1.

¹Cf. M. Weinhardt and W. Luk: *Pipeline Vectorization*, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, Feb. 2001, pp. 234-248.

3 Language Coverage

This Section describes which C files can currently be handled by XPP-VC.

3.1 Restrictions

3.1.1 XPP Restrictions

The following C language operations cannot be mapped to an XPP Core at all. They are not allowed in XPP-VC programs and need to be mapped to the host processor in a codesign compiler, cf. Section 7.3,

- Operating system calls, including I/O
- Division, modulo, non-constant shift and floating point operations (unless XPP Core's ALU supports them)²
- The size of arrays mapped to internal RAMs is limited by the number and size of internal RAM blocks.

3.1.2 XPP-VC Compiler Restrictions

The current XPP-VC implementation necessitates the following restrictions:

1. No multi-dimensional constant arrays (due to the SUIF version currently used)
2. No switch/case statements
3. No struct data types
4. No function calls except the XPP port and pragma functions defined in Section 3.2.1. The program must only have one function (main).
5. No pointer operations
6. No library calls or recursive calls
7. No irregular control flow (break, continue, goto, label)

Additionally, there are currently some implementation-dependent restrictions for vectorized loops, cf. the Release Notes. The compiler produces an explanatory message if an inner loop cannot be pipelined despite the absence of dependences. However, for many of these cases, simple workarounds by minor program changes are available. Furthermore, programs which are too large for one configuration cannot be handled. They should be split into several configurations and sequenced onto the XPP Core, using NML's reconfiguration commands. This will be performed automatically in later releases by temporal partitioning, cf. Section 7.1.

²In future XPP-VC releases, an alternative, sequential implementation of these operations by NML macros will be available.

3.2 XPP-VC C Language Extensions

We now describe C language extensions used by XPP-VC. In order to use these extensions, the C program must contain the following line:

```
#include "XPP.h"
```

This header file, *XPP.h*, defines the port functions defined below as well as the pragma function `XPP_unroll()`. If `XPP_unroll()` directly precedes a FOR loop, it will be completely unrolled by *partition*, cf. Section 6.2.

3.2.1 XPP Port Functions

Since the normal C I/O functions cannot be used on an XPP Core, a method to access the XPP I/O units in port mode is provided. *XPP.h* contains the definition of the following two functions:

```
XPP_getstream(int ionum, int portnum, int *value)
```

```
XPP_putstream(int ionum, int portnum, int value)
```

ionum refers to an I/O unit (1..4), and *portnum* to the port used in this I/O unit (0 or 1). For the duration of the execution of a program, an I/O unit may only be used either for port accesses or for RAM accesses (see below). If an I/O unit is used in port mode, each *portnum* can only be used either for read or for write accesses during the entire program execution. In the access functions, *value* is the data received from or written to the stream. Note that `XPP_getstream` can currently only read values into scalar variables (not directly into array elements!), whereas `XPP_putstream` can handle any expressions. An example program using these functions is presented in Section 6.1.

3.2.2 pragma Directives

Arrays can be allocated to external memory by a compiler directive:

```
#pragma extern <var> <RAM_number>
```

Example: `#pragma extern x 1` maps array *x* to external memory bank 1.

Note the following:

- *<var>* must be defined *before* it is used in the pragma.
- Bank *<RAM_number>* must be declared in the file *xppvc_options*, cf. Section 4.
- If two arrays are allocated to the same external RAM bank, they are arranged in the order of appearance of their respective pragma directives. The resulting offsets are recorded in *file.isf*, cf. Section 5.1.

4 Directories and Files

After correct installation, the `XPPC.ROOT` environment variable is defined, and the `PATH` variable extended. `$XPPC.ROOT` is the XPP-VC root directory. `$XPPC.ROOT/bin` contains all binary files and the scripts `xppvcmake` and `xppgcc`. `$XPPC.ROOT/doc` contains this manual and the file `xppvc_releasenotes.txt`. `XPP.h` is located in the include subdirectory.

Finally, `$XPPC.ROOT/lib` contains the options file `xppvc_options`. If an options file with the same name exist in the current working directory or the `.xds` subdirectory of the user's home directory, they are used (in this order) instead of the master file in `$XPPC.ROOT/lib`.

Option	Explanation	Default value in <code>xppvc_options</code>
<code>debug</code>	debug output enabled	on
<code>version</code>	XPP IP version	V2
<code>pacsize</code>	number of ALU-PAEs in x and y direction	6 / 12
<code>xppsiz</code>	number of PACs in x and y direction	1 / 1
<code>busnumber</code>	number of data and event buses per row (both dir.s)	6 / 6
<code>iramsize</code>	number of words in one internal RAM	256
<code>bitwidth</code>	XPP data bit width	32
<code>freg_data_port</code>	number of FREG data ports	3
<code>breg_data_port</code>	number of BREG data ports	3
<code>freg_event_port</code>	number of FREG event ports	4
<code>breg_event_port</code>	number of BREG event ports	4

Table 1: Options

`xppvc_options` sets the compiler options listed in Table 1. Most of them define the XPP IP parameters which are used in the generated NML file. Lines starting with a `#` character are comment lines.

Additionally, `extram` followed by four integers declares the external RAM banks used for storing arrays. At most four external RAMs can be used. Each integer represents the size of the bank declared. Size zero must be used for banks which do not exist. The master file contains the following line which declares four 4GB (1 G words) external banks:

```
extram 1073741824 1073741824 1073741824 1073741824
```

Note that, in order to simplify programming, `xppvc_options` does not have to be changed if an I/O unit is used for port accesses. However, this memory bank is not available in this case despite being declared.

5 Using XPP-VC

5.1 xppvcmake

In order to create an NML file, `file.c` is compiled with the command `xppvcmake file.nml`. `xppvcmake file.xbin` additionally calls `xmap`. With `xppvcmake`, `XPP.h` is automatically searched

for in directory \$XPPC.ROOT/include.

The following output produced by translating the example program *streamfir.c* in Section 6.1 shows the programs called by *xppvcmake*:

```
$ xppvcmake streamfir.nml
pscc -I/home/wema/xppc/include -parallel -no PORKY_FORWARD_PROP4
-spr streamfir.c
porky -dead-code streamfir.spr streamfir.spr2
partition streamfir.spr2 streamfir.svo
Program analysis:
  main: DO-LOOP, line 9 can be synthesized
  main: can be synthesized completely
Program partitioning:
  Entire program selected for XPU module synthesis.
  main: DO-LOOP, line 9 selected for synthesis
porky -const-prop -scalarize -copy-prop -dead-code streamfir.svo
streamfir.svol
predep -normalize streamfir.svol streamfir.svo2
porky -ivar -know-bounds -fold streamfir.svo2 streamfir.sur
nmlgen streamfir.sur streamfir.xco
```

pscc is the SUIF frontend which translates *streamfir.c* into the SUIF intermediate representation, and *porky* performs some standard optimizations. Next, *partition* analyzes the program. The output indicates that the entire program can and will be mapped to NML. Then *porky* and *predep* perform some additional optimizations before *nmlgen* actually generates the file *streamfir.nml*. The SUIF file *streamfir.xco* is generated to inspect and debug the result of code transformations.³ In the generated NML file, only the I/O ports are placed. All other objects are placed automatically by *xmap*. Cf. Section 6.1 for an example of the *xsim* program using the I/O ports corresponding to the stream functions used in the program.

For an input file *file.c*, *nmlgen* also creates an interface description file *file.itf* in the working directory. It shows the array to RAM mapping chosen by the compiler. In the *debug* subdirectory (which is created), files *file.part.dbg* and *file.nmlgen.dbg* are generated. They contain more detailed debugging information created by *partition* and *nmlgen* respectively. The files *file.first.dot* and *file.final.dot* created in the *debug* directory can be viewed with the *dotty* graph layout tool. They contain graphical representations of the original and the transformed and optimized version of the generated control/dataflow graph.

5.2 xppgcc

This command is provided for comparing simulation results obtained with *xppvcmake*, *xmap* and *xsim* (or from execution on actual XPP hardware) with a "direct" compilation of the C program with *gcc* on the host. *xppgcc* compiles the input program with *gcc* and binds it with predefined *XPP_getstream* and *XPP_putstream* functions. They read or write files *port<n>-<m>.dat* in the

³In an extended codesign compiler, the *.xco* file would also be used to generate the host partition of the program.

current directory for n in 1..4 and m in 0..1. For instance, the program in Section 6.1 is compiled as follows:

```
xppgcc -o streamfir streamfir.c
```

The resulting program *streamfir* will read input data from *port1_0.dat* and write its results to *port4_0.dat*.⁴

6 Examples

6.1 Stream Access

The following program *streamfir.c* is a small example showing the usage of the `XPP_getstream` and `XPP_putstream` functions. The infinite `WHILE`-loop implements a small FIR filter which reads input values from port 1_0 and writes output values to port 4_0. The variables `xd`, `xdd` and `xddd` are used to store delayed input values. The compiler automatically generates a shift-register-like configuration for these variables. Since no operator dependences exist in the loop, the loop iterations overlap automatically, leading to a pipelined FIR filter execution.

```
1 #include "XPP.h"
2
3 main() {
4 int x, xd, xdd, xddd;
5
6 x = 0;
7 xd = 0;
8 xdd = 0;
9 while (1) {
10 xddd = xdd;
11 xdd = xd;
12 xd = x;
13 XPP_getstream(1, 0, &x);
14 XPP_putstream(4, 0, (2*x + 6*xd + 6*xdd + 2*xddd) >> 4);
15 }
16 }
```

After generating *streamfir.xbin* with the command `xppvcmake streamfir.xbin`, the following command reads the input file *port1_0.dat* and writes the simulation results to *xpp_port4_0.dat*.

```
xsim -run 2000 -in1_0 port1_0.dat -out4_0 xpp_port4_0.dat
streamfir.xbin > /dev/null
```

xpp_port4_0.dat can now be compared with *port4_0.dat* generated by compiling the program with `xppgcc` and running it with the same *port1_0.dat*.

⁴However, programs receiving initial data from or writing result data to external RAMs in `xsim` cannot be compared to directly compiled programs using `xppgcc`. The results may also differ if a bitwidth other than 32 is used for the generated NML files.

6.2 Array Access

The following program *arrayfir.c* is an FIR filter operating on arrays. The first FOR-loop reads input data from port 1_0 into array *x*, the second loop filters *x* and writes the filtered data into array *y*, and the third loop outputs *y* on port 4_0.

```

1 #include "XPP.h"
2 #define N 256
3 int x[N], y[N];
4 const int c[4] = { 2, 4, 4, 2 };
5 main() {
6     int i, j, tmp;
7     for (i = 0; i < N; i++) {
8         XPP_getstream(1, 0, &tmp);
9         x[i] = tmp;
10    }
11    for (i = 0; i < N-3; i++) {
12        tmp = 0;
13        XPP_unroll();
14        for (j = 0; j < 4; j++) {
15            tmp += c[j]*x[i+3-j];
16        }
17        y[i+2] = tmp;
18    }
19    for (i = 0; i < N-3; i++)
20        XPP_putstream(4, 0, y[i+2]);
21 }

```

xppvcmake produces the following output:

```

$ xppvcmake arrayfir.nml
pscc -T/home/wema/xppc/include -parallel -no PORKY FORWARD PROP4
-.spr arrayfir.c
porky -dead-code arrayfir.spr arrayfir.spr2
partition arrayfir.spr2 arrayfir.svo
Program analysis:
  main: FOR-LOOP i, line 7 can be synthesized/vectorized
  main: FOR-LOOP j, line 14 can be synthesized/unrolled/vectorized
  main: FOR-LOOP i, line 11 can be synthesized/vectorized
  main: FOR-LOOP i, line 19 can be synthesized/vectorized
  main: can be synthesized completely
Program partitioning:
  Entire program selected for NML module synthesis.
  main: FOR-LOOP i, line 7 selected for pipeline synthesis
  main: FOR-LOOP i, line 11 selected for pipeline synthesis
  main: FOR-LOOP i, line 19 selected for pipeline synthesis

```

Vectorizing C Compiler XPP-VC

9

```

... unrolling loop j
porky -const-prop -scalarize -copy-prop -dead-code array-
fir.svo arrayfir.svo1
predep -normalize arrayfir.svo1 arrayfir.svo2
porky -ivar -know-bounds -fold arrayfir.svo2 arrayfir.sur
nmigen arrayfir.sur arrayfir.xcu

```

The messages from partition show that all loops can be vectorized. The dependence analysis did not find any loop-carried dependences preventing vectorization. The inner loop in the middle of the program is unrolled. The outer loop's body is effectively substituted by the following statement:

```

y[i+2] = c[0]*x[i+3] + c[1]*x[i+2] + c[2]*x[i+1] + c[3]*x[i];

```

Since all remaining loops are innermost loops, they are selected for pipeline synthesis. Array reads, computations, and array writes overlap. To reduce the number of array accesses, the compiler automatically removes redundant array reads. In the middle loop, only $x[i+3]$ is read. For $x[i+2]$, $x[i+1]$ and $x[i]$, delayed versions of $x[i+3]$ are used, forming a shift-register. Therefore, each loop iteration needs only one cycle since one read from x , all computations, and one write to y can be executed concurrently.

Finally, the following example program fragment is a 2-D edge detection algorithm.

```

/* 3x3 horiz. + vert. edge detection in both directions */
for(v=0; v<=VERLEN-3; v++) {
  for(h=0; h<=HORLEN-3; h++) {
    htmp = (p1[v+2][h] - p1[v][h]) +
            (p1[v+2][h+2] - p1[v][h+2]) +
            2 * (p1[v+2][h+1] - p1[v][h+1]);
    if (htmp < 0)
      htmp = - htmp;

    vtmp = (p1[v][h+2] - p1[v][h]) +
            (p1[v+2][h+2] - p1[v+2][h]) +
            2 * (p1[v+1][h+2] - p1[v+1][h]);
    if (vtmp < 0)
      vtmp = - vtmp;

    sum = htmp + vtmp;
    if (sum > 255)
      sum = 255;
    p2[v+1][h+1] = sum;
  }
}

```

As the output of partition shows, both loops can be vectorized. Since only innermost loops can be pipelined, the outer loop is executed sequentially. (Note that the line numbers in the program outputs are not obvious since only a program fragment is shown above.)

Vectorizing C Compiler XPP-VC

partition edge.spr2 edge.svo

Program analysis:

main: FOR-LOOP h, line 22 can be synthesized/can be vectorized

main: FOR-LOOP v, line 21 can be synthesized/can be vectorized

main: can be synthesized completely

Program partitioning:

Entire program selected for XPP module synthesis.

main: FOR-LOOP h, line 22 selected for pipeline synthesis

main: FOR-LOOP v, line 21 selected for synthesis

Also note the following additional features of this program: Address generators for the 2-D array accesses are automatically generated, and the array accesses are reduced by generating shift-registers for each of the three image lines accessed. Furthermore, the conditional statements are implemented using SWAP (MUX) operators. Thus the streaming of the pipeline is not affected by which branch the conditional statements take.

7 Future Compiler Extensions

Apart from removing some of the restrictions of Section 3.1.2, the following extensions are planned for XPP-VC.

7.1 Temporal Partitioning

By using the pragma function *XPP_next_conf()*, programs are partitioned into several configurations which are loaded and executed sequentially on the XPP Core. Specific NML configuration commands are generated which also exploit XPP's sophisticated configuration and preloading capabilities. Eventually, the temporal partitions will be determined automatically.

7.2 Program Transformations

For more efficient XPP configuration generation, some program transformations are useful. In addition to loop unrolling, *loop merging*, *loop distribution* and *loop tiling* will be used to improve loop handling, i.e. enable more parallelism or better XPP usage.

Furthermore, programs containing more than one function could be handled by *inlining function calls*.

7.3 Codesign Compiler

This section sketches what an extended C compiler for an architecture consisting of an XPP Core combined with a host processor might look like. The compiler should map suitable program parts, especially inner loops, to the XPP Core, and the rest of the program to the host processor. I.e., it is a host/XPP codesign compiler, and the XPP Core acts as a coprocessor to the host processor.

This compiler's input language is full standard ANSI C. The user uses pragmas to annotate those program parts that should be executed by the XPP Core (manual partitioning). The compiler checks if the

Public

2002-1-11 ccomp V1.1

Vectorizing C Compiler XPP-VC

11

selected parts can be implemented on the XPP. Program parts containing non-mappable operations must be executed by the host.

The program parts running on the host processor ("SW"), and the parts running on the PAE array ("XPP") cooperate using predefined routines (*copy_data_to_XPP*, *copy_data_to_host*, *start_config(n)*, *wait_for_coprocessor_finish(n)*, *request_config(n)*). For all XPP program parts, XPP configurations are generated. In the program code, the XPP part *n* is replaced by *request_config(n)*, *start_config(n)*, *wait_for_coprocessor_finish(n)*, and the necessary data movements. Since the SUIF compiler contains a C backend, the altered program (host parts with coprocessor calls) can simply be written back to a C file and then processed by the native C compiler of the host processor.

Thus the sequential control flow of the C program defines when XPP parts are configured into the XPP Core and executed.

Fast and Guaranteed C-Compilation onto the PACT-XPP Reconfigurable Computing Platform

João M. P. Cardoso, and Markus Weinhardt

PACT Informationstechnologie GmbH
Leopoldstr. 236, D-80807 München, Germany
(cajo, mw)@pactcorp.com

Abstract

The eXtreme Processing Platform (XPP) technology offers a unique reconfigurable computing platform supported with a set of tools. A C compiler, which integrates both new and efficient compilation techniques and temporal partitioning, is presented. Temporal partitioning guarantees the compilation of programs with unlimited complexity as long as the supported C-subset is used. A new partitioning scheme, which permits to map large loops of any kind and is neither constrained by loop-dependencies nor nested structures, is also presented. Furthermore, temporal partitioning is applied to reduce the configuration time overhead and thus can lead to performance gains. The compilation from C code to the configuration data, ready to be downloaded onto the XPP, takes seconds for complex examples, which is, as far as we know, not reproduced by any other reconfigurable computing technology. The compiler represents a step forward, by furnishing a truly "push-button" approach, only comparable to microprocessor domains, and thus can spread the use of the XPP technology and deal with time-to-market pressures positively.

1 Introduction

Many of today's applications are characterized by intensive data-stream processing and high-performance requirements. Such performance is more and more evident to not be accomplished with today's microprocessor technology. Conventional processors (including DSPs) are geared for sequential processing. Multi-DSP and very large instruction word (VLIW) processors still have severe memory bottlenecks, lack the number of data ports required to support multi-channel, high speed data streams, and fail on furnishing low power consumption solutions. Accelerating specific functions using application-specific integrated circuits (ASICs) relieves some of the processing burden, adds some required features, but limits flexibility and requires expensive non-recurring engineering (NRE) costs and long design

cycles. High density field-programmable gate arrays (FPGAs) eliminate the NRE costs, add flexibility, but still require long timing optimizations and verification cycles and low level hardware efforts. Additionally, the fine-grained structure adopted in FPGAs is not suitable to map at the algorithmic level, which is proved by the well-known difficulties to have a "push-button" high-level methodology to program these architectures.

New reconfigurable processing units (RPU)s are being introduced trying to solve those problems [1]. One of the new promising architectures is the XPP [2][3]. The XPP is a coarse-grained, runtime-reconfigurable, 2-D array parallel structure. The architecture was designed to facilitate programming and to support pipelining, dataflow computations, and parallelism from the instruction to the task level efficiently. Therefore, this technology is well suited for applications in multimedia, telecommunications, simulation, digital signal processing, and similar stream-based application domains. The XPP architecture also supports dynamic self-reconfiguration in a user transparent way. In order to drastically reduce the time to program the XPP, and to keep the user from architecture details, a high-level compiler integrating temporal partitioning is required. Such a compiler is the main topic of this paper.

This paper is organized as follows. The next section introduces briefly the XPP technology. Section 3 outlines compilation to the XPP and section 4 describes the temporal partitioning steps. Section 5 shows some experimental results, section 6 points out the main differences between this and previous works, and finally section 7 concludes the paper and enumerates ongoing and future work planned.

2 XPP Technology

The XPP technology consists of a reconfigurable computing platform delivered as a device or an intellectual property (IP) core, and a complete development tool suite (XDS) [2]. An XPP can be used as a coprocessor for CPU and DSP

architectures. A prior version of the technology has resulted in the XPU128-ES [3], a prototype device, which was produced in silicon.

The XPP architecture is based on a hierarchical array of coarse-grain, adaptive computing elements called *Processing Array Elements (PAEs)*, and a *packet-oriented communication network*. The strength of the XPP technology originates from the combination of array processing with unique and powerful run-time reconfiguration mechanisms. Different tasks or applications can be configured and run independently on different parts of the array. Reconfiguration is triggered externally or even by special event signals originating within the array, enabling self-reconfiguring designs. By utilizing protocols implemented in hardware, data and event packets are used to process, generate, decompose and merge streams of data.

2.1 Array Structure

An XPP contains one or several *Processing Array Clusters (PACs)*, i.e., rectangular blocks of PAEs. Fig.1 shows the structure of a typical XPP device. It contains four PACs (see top left hand side). Each PAC is attached to a *Configuration Manager (CM)* responsible for writing configuration data into the configurable objects of the PAC using a dedicated bus. Multi-PAC XPPs contain additional CMs for configuration data handling, forming a hierarchical tree of CMs. The root CM is called the supervising CM (SCM). It has an external interface (dotted arrow originating from the SCM in Fig.1) which usually connects the SCM to an external configuration memory. A CM consists of a state machine and internal RAM for configuration caching (see top right-hand side of Fig.1).

Horizontal busses carry data and events. They can be segmented by configurable switch-objects, and connected to PAEs and special I/O objects at the periphery of the device. The I/O objects can be used for data-streaming or to access external resources (e.g., memories). A column of ports to the corresponding leaf CM is located on the array. A CMPort can be used to send events to the CM from the array. The typical PAE shown in Fig.1 (bottom center) contains three objects: one FREG (forward register), one BREG (backward register) and one ALU. The FREG object is used for vertical forward routing (with a programmable number of register stages), or to perform MERGE, SWAP or DEMUX operations (for controlled stream manipulations). The BREG object is used for vertical backward routing (registered or not), or to perform some selected arithmetic operations (e.g., ADD, SUB, SHIFT). The DREGs can also be used to perform logical operations on events. Each ALU (see its internal structure on the bottom left-hand side of Fig.1) performs common two-input fixed-point arithmetical and logical operations, and com-

parisons. A MAC (multiply and accumulate) operation can be performed using the ALU and the BREG objects of one PAE in a single clock cycle.

Another standard PAE object is the memory object which can be used in FIFO mode or as RAM for lookup tables, intermediate results, etc. If such objects are needed they are located in the left and/or right columns of PAEs of each PAC. However, any PAE object functionality can be included in the XPP architecture.

A set of parameterizable features can be used to furnish an XPP that best fits to user and application demands. Those features include: the number of PACs and their PAEs, number of internal memories, number of I/O ports, number of buses, word bitwidth, cache size, depth of the FIFO to configure each object, etc.

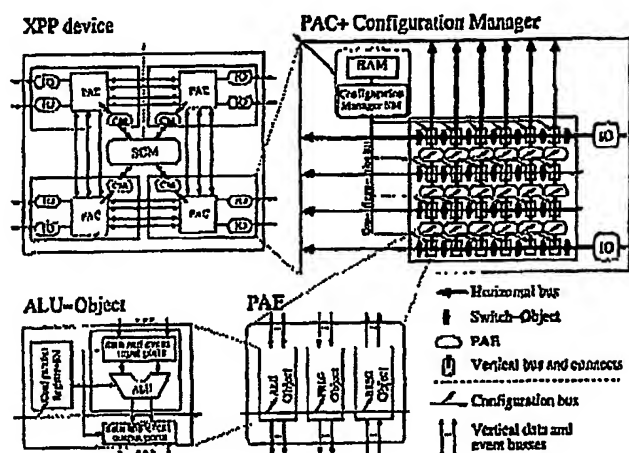


Figure 1: XPP architecture.

2.2 Packet Handling and Synchronization

PAE objects as defined above communicate via a packet-oriented network. Two types of packets are sent through the array: data and event packets. Data packets have a uniform bitwidth specific to the XPP Core or device.

In normal operation mode, PAE objects are self-synchronizing. An operation is performed as soon as all necessary data input packets are available. The results are forwarded as soon as they are computed and the previous results have been consumed. Thus, a signal-flow graph can be mapped directly to the ALU objects and data-streams can flow through them in a pipelined manner without adding specific hardware.

Event packets are one bit wide. They transmit state information which controls ALU execution and packet generation. For instance, they can be used to control the merg-

ing of data-streams or to deliberately discard data packets. Thus, conditional computations depending on the results of earlier ALU operations are feasible. Events can even trigger a self-reconfiguration of the device as explained below.

Each data or event packet is only forwarded if the previous one has already been consumed. The communication system was designed to transmit one packet on each interconnect per cycle. Hardware protocols ensure that no packets are lost, even in the case of pipeline stalls or during the configuration process. This simplifies application development considerably. No explicit scheduling of operations is required.

2.3 Configuration

The XPP architecture is optimized for rapid and user-transparent configuration. For this purpose, the configuration managers in the CM tree operate independently (without global synchronization), and therefore are able to configure their respective parts of the array in parallel. Every PAE stores locally its current configuration state, i.e., if it is part of a configuration or not (states "configured" or "free"). Once a PAE is configured, it changes its state to "configured". This prevents the respective CM from reconfiguring a PAE which is still in use. The CM caches the configuration data in its internal RAM and constantly tries to configure the objects used by the next configuration requested. Each XPP object has a configuration FIFO which stores data of subsequent configurations. Once an object has been released (state "free"), the next configuration word in its FIFO is loaded immediately. Hence it is possible to reconfigure partially in one clock cycle. Additionally, a prefetching mechanism is used. While a configuration is being loaded onto the FIFO of each object, other configurations may already be requested and cached in the low-level CMs' internal RAM. Thus, it does not need to be requested all the way from the SCM down to the array when objects become available. While loading a configuration, its PAEs start their part of the computations as soon as they are in state "configured".

Each ALU object has an input event port that triggers the self-releasing of its resources and of all of the objects connected to it. Such event is successively broadcasted according to the interconnections.

Because of its coarse-grain nature, an XPP device can be configured rapidly. Since only the configuration of those array objects actually used is necessary, the configuration time depends on the application.

2.4 Development Tools

The XPP can be programmed by using the *Native Mapping Language* (NML) [2], a PACT proprietary structural

language with reconfiguration primitives. It gives the programmer direct access to all hardware features. In NML, configurations consist of modules which are specified as in a structural hardware description language, similar to, for instance, structural VHDL. PAE objects are explicitly allocated, optionally placed, and their connections specified. Additionally, NML includes statements to support configuration handling. Thus, configuration handling is an explicit part of the NML application program. XDS is an integrated environment for programming with NML. The main component is the mapper *xmap* which compiles NML source files, places and routes the objects, and generates XPP binary files. *xmap* uses an enhanced force-based placer with short runtimes. The XPP binaries can either be simulated and visualized cycle by cycle with the *xsim* and *xvis* tools, or directly executed on an XPP device. A high-level compiler, described in the next section, has been added to XDS and permits to map C programs onto the XPP.

2.5 Application Execution on XPP

Reconfiguration and prefetching requests can be issued by any CM in the tree (including the SCM which can respond to external requests) and also by event signals generated in the array itself. Running modules can do a self-releasing of their resources and request another configuration. Thus, it is possible to execute an application consisting of several configurations without any external control.

The CM of the XPP permits to exploit speculative configuration¹, i.e., the configuration of a module possibly used after the current one has finished execution. If the path which includes that module is taken, the CM only has to trigger the execution of the configuration (see the section of the NML code in Fig.2 and the simulation performed with *xsim* in Fig.3, where *conf_MOD2* is speculatively configured during the execution of *conf_MOD0*). If this path is not taken, the CM triggers the releasing of the resources already configured and requests the other configuration.

3 Compiling C Code with XPP-VC

The XPP Vectorizing C Compiler XPP-VC is based on the SUIF compiler framework [4]. SUIF is used because of its easily extensible properties. The XPP-VC compilation flow is shown in Fig.4. An options file, used by the compiler, specifies the parameters of the targeted XPP and the external memories connected to the XPP. To access XPP I/O ports specific C-functions are provided.

¹This has similarities to speculative execution. In this case, before knowing if a configuration will be requested, its configuration is started.

```

...
CONFIG conf_MOD0 {
  CONF_MODULE(MOD0) // request the configuration of MOD0
  REQUEST(conf_MOD2_spec) // start speculative configuration
  // if (MOD0.CMPort0 == "0") then conf_MOD2_exec is requested
  CONF_CMPORT(MOD0.CMPort0, conf_MOD2_exec, _)
  // if (MOD0.CMPort1 == "0") then conf_MOD1 is requested
  CONF_CMPORT(MOD0.CMPort1, conf_MOD1, _)
}
CONFIG conf_MOD2_spec { // request the configuration of MOD2
  CONF_MODULE(MOD2) // but do not start it
}
CONFIG conf_MOD2_exec { // MOD2 is taken
  SET(MOD2.StartA = 1) // enable the start of computing of MOD2
  REQUEST(conf_MOD3) // request the next configuration
}
CONFIG conf_MOD1 { // MOD1 is taken
  REQUEST(conf_MOD2_rec) // releasing of resources of MOD2
  CONF_MODULE(MOD1) // request the MOD1
  REQUEST(conf_MOD3) // request the next configuration
}
CONFIG conf_MOD2_rec {
  RECONF(MOD1.Start) // release the resources of MOD1
}
...

```

Figure 2: Section of NML describing the control flow.

The compiler starts with some architecture-independent preprocessing passes based on well-known compilation techniques [5]. During this step, FOR loops are automatically unrolled if instructed by the programmer. Then the compiler performs a data-dependence analysis. The compiler tries to vectorize inner program FOR-loops. In XPP-VC, vectorization means that loop iterations are overlapped and executed in a pipelined, parallel fashion. This technique is based on the *Pipeline Vectorization* method developed for reconfigurable architectures [6].

The C program can be manually splitted in several modules by using annotations. Otherwise, automatic temporal partitioning can be applied (see section 4) in order to furnish mappable modules and to reduce the overall latency.

MODGen generates one NML module for each temporal partition. First, program data is allocated on the XPP. By default, MODGen maps each program array to internal or external RAM while scalar variables are stored in registers within the PAEs. Next, a control/dataflow graph (CDFG) is generated. Straight-line code without array accesses can be directly mapped to a data-flow graph since the data dependences are obvious in the DAG representation. One ALU is allocated for each operator in the CDFG. Because of the self-synchronization of operators on the XPP, no explicit control or scheduling is needed. The same is true for conditional execution of such blocks. Both branches are executed in parallel and MUX operators select the correct output (and discard the other one) depending on the condition. This data-driven execution of the operators automatically yields instruction-level parallelism. In contrast, accesses to the

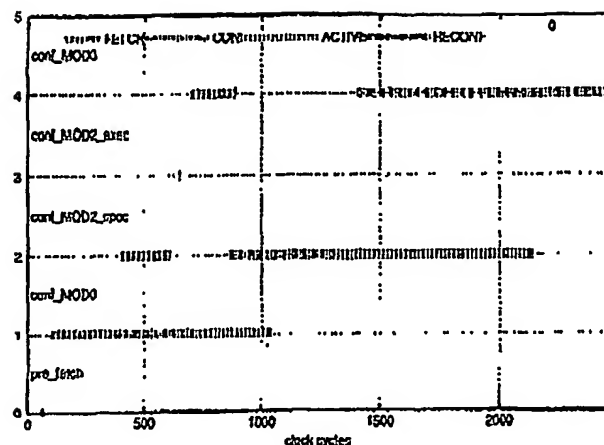


Figure 3: Speculative configuration (enables earlier activation of MOD2).

same array have to be controlled explicitly to maintain the correct execution order. MERGE operators (which select one input without discarding the other one) route address and write data packets in the correct order to the RAM, and DEMUX operators route read data packets to the correct subsequent operator. State machines for generating the correct sequence of event signals (to control these operators) are synthesized by the compiler. For conditional branches, containing array accesses or inner loops, DEMUX operators controlled by the IF condition route data packets only to the selected branch, and output values are taken from the branch activated. Thus, only selected branches receive data packets and execute.

In loops, all variables updated in the loop body are handled as follows. The first iteration uses an input packet for the variable's value, and the subsequent iterations use packets generated in the previous iteration. In all but the last iteration, a DEMUX operator routes the outputs of the loop body back to the body inputs. Only the results of the last iteration are routed to the loop output by the DEMUX operators. The control packets for the DEMUX are generated by the loop counter or the comparator evaluating the exit condition. Note that the internal operators' outputs cannot just be connected to subsequent operators since they produce a result in each loop iteration. The required last packet would be hidden by a stream of intermediate packets. If array accesses are present, a loop iteration may only be started after the previous iteration has terminated because the original access order must be maintained. This is enforced by event signals.

For generating more efficient XPP configurations, MODGen generates pipelined operator networks for inner pro-

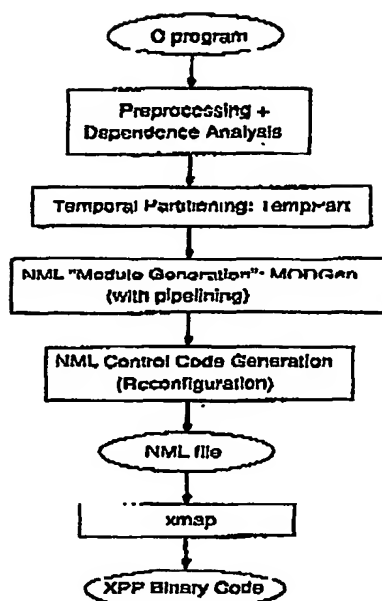


Figure 4: XPP-VC compilation flow.

grum loops which have been annotated for vectorization by the preprocessing step. In other words, subsequent loop iterations are started before previous iterations have finished. Data packets flow continuously through the operator pipelines. By applying pipeline balancing techniques, maximum throughput is achieved. For many programs, additional performance gains are achieved by the complete loop unrolling transformation. Although unrolled loops require usually more XPP resources, they yield more parallelism and better exploitation of the XPP. To reduce the number of array accesses, the compiler automatically removes redundant array reads. When array references inside loops access subsequent element positions the compiler only uses one reference and generates delayed structures, forming shift-registers.

Finally, each module generated by MODGen is placed and routed automatically by xmap.

The XPP-VC compiler currently supports a C-subset sufficient for programming real applications. Struct data types, pointer operations, irregular control flow (break, continue, goto, label), and recursive and operating system calls are not supported or cannot be mapped to the XPP.

4 Temporal Partitioning

A program too large to fit in an XPP can be handled by splitting it in several parts (configurations) such that each one is mappable. Temporal partitioning permits the automatic exposing of configurations such that the overall execution time of the application is minimized and is successfully mapped onto the XPP resources. It considers the costs to load into the cache, to configure and to execute each configuration with the XPP. An important strategy that is considered is to pre-fetch configurations while another is being configured or is running. Arrays of constants or with pre-defined values used in one or more configurations can be initialized in parallel with the execution of the previous configurations. This takes advantage of the initialization of the array carried out by using the configuration bus.

The set of partitions resulting from the splitting are then processed by MODGen, generating a set of configurations. Next, specific NML configuration commands are generated which also exploit XPP's sophisticated configuration and pre-fetching capabilities, and specify the configuration control flow that is orchestrated by the CM.

4.1 Benefits of Temporal Partitioning

Temporal partitioning targeting the XPP can reduce, when efficiently applied, the overall execution time. Such reduction can be mainly achieved by the following issues: (1) reduction of each partition complexity can reduce the interconnection delays (long interconnections may pass through registers and thus add clock cycle delays); (2) reduction of the number of references, in the section of the program related to each partition, using the same resource, by distributing the overall references among partitions, can lead to performance gains as well. This happens with the statements presented in the program referring the same array; (3) reduction of the overall configuration overhead by overlapping fetching, configuration and execution of distinct partitions.

Example: Consider the C example `max_avg` shown in Fig.5. Configuration boundaries are represented by `XPP_next_conf()` statements. They define four configurations in the code (see Fig.6). Apart from exposing temporal partitions in such a way that the mapping to XPP is accomplished, combining only the most frequently taken conditional paths in the same partition can reduce the total execution time by substantially reducing the reconfiguration time (since the partitions for the other paths are not configured when they are not taken). Fig.6 presents such a case. If the path `bb_0` and `bb_1` has been identified as the most frequently executed, this path can be in the same partition².

²Tail duplication of `bb_3` would permit to have a configuration with (`bb_0`, `bb_1`, `bb_3`) and another one with (`bb_2`, `bb_3`);

In such a case, the configuration related to `bb_2` will only be called when the most frequent path has not been taken.

```
// max_avg example
...
if(op==1) { // average kernel
    sum = 0;
    for(i=0; i<N; i++) {
        sum+=x[i];
    }
    average = sum/N;
    XPP_next_conf();
} else { // max kernel
    XPP_next_conf();
    max = 0;
    for(i=0; i<N; i++) {
        if(x[i] > max) max = x[i];
    }
    XPP_next_conf();
}
...
```

Figure 5: Example with two conditionally executed kernels and with configuration boundaries represented.

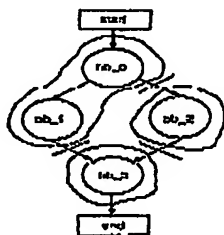


Figure 6: CFG of the algorithm shown in Fig. 5. Lines crossing edges represent the `XPP_next_conf()` statements in the code. Bubbles containing basic blocks represent the regions to be implemented in different partitions.

Since configuration takes many clock cycles, it is in most cases preferable to reuse a configuration as much long as possible in order to reduce the reconfiguration time overhead. Thus, loops in the source code are always good candidates to be entirely implemented by a single configuration.

4.2 Partitioning Loops

Each loop that does not fit onto the XPP can be dealt with by performing loop distribution [5] (if applicable) or by partitioning the loop and use the CM to orchestrate the control flow. Currently, loop distribution is not automatically applied. Instead, we propose a new method to partition complex loops without resolutions. All the loops which their bodies must be partitioned are transformed into straight line code with a jump to loop-exit or to the next iteration in order that each partition can be compiled by MODGen. Fig. 7

shows an example of such transformation without the statements needed to communicate the value of scalar variables between configurations. Each configuration requests the next configuration to be taken (if none is requested then the application terminates and the last configuration releases its resources). Depending on the value of the `i<N` condition, config. #2 takes two different exits, which requests #3 or #4 respectively. Since config. #3 always requests #2, at the end of its execution, the initial behavior of the loop is preserved. The temporal partitioning creates two additional configuration boundaries to preserve the initial functionality. From Fig 7b can be seen that configuration boundaries were inserted before and after the `if` statement. These boundaries are needed since the code before and after will be executed once and both the `if` header and body will iterate `N+1` and `N` times respectively.

int i;	int i;	
...	i=0;	#1
for(i=0; i<N; i++) {	lab1: if(i<N) {	#2
stmt1;	stmt1;	#2
XPP_next_conf();	stmt2;	#3
stmt2;	i++;	#3
}	goto lab1;	#3
exit3;	}	#3
a)	stmt3;	#4
	b)	c)

Figure 7: Example of the transformation applied for partitioning loops. a) original code added with the statement representing where the loop is partitioned; b) transformed code; c) configuration ID for each statement in b).

4.3 Automatic Partitioning

From the SUIF representation of the C source code the temporal partitioning phase constructs an Hierarchical Task Graph³ extended, HTG+. This extended graph has two types of nodes: (1) *behavioral nodes* representing lines of code in the input program; (2) *array nodes* representing each array existent in the source code. For instance, Fig. 8 shows the top level of the HTG+ for an implementation of the DCT (Discrete Cosine Transform) based on matrix multiplications. Type (1) nodes have three distinct sub-types: (a) *block nodes* representing basic blocks; (b) *compound nodes* representing `if-then-else` structures; (c) *loop nodes* representing the loops (`for`, `while`). Loop and compound nodes explicitly embody hierarchical levels. Edges in the HTG+ represent data communication between two nodes or just enforce execution's precedence.

Each behavioral node of the HTG+ is labeled with the following information (some of the labelling steps require

³The model has been chosen, because it also exposes loop and task level parallelism.

estimation efforts): (1) *block and compound nodes*: number of ALUs and REGs; (2) *loop nodes*: number of iterations (if unbound, profiling can be used), and number of ALUs and REGs; (3) *array nodes*: the size of the array, type of the elements, and, when they do exist, the initialization values. Each edge between two behavioral nodes of the HTG+ is labeled with the number of data words that must be transferred between the two nodes. Each edge between an array and a behavioral node in the HTG+ is labeled with the number of load and store references in the source code represented by the behavioral node to that particular array. The estimated number of times that each load and store reference will be executed is also collected. The use of the same array by different behavioral nodes, increases the execution latency and the number of resources needed for this partition⁴.

TempPart uses three types of estimations: (1) number of XPP resource units needed by the configuration implementing a single or a set of behavior nodes; (2) latency for a behavior node or a set of connected behavior nodes on the HTG+ (this does not need to be accurate to the real execution time and only needs to have relative accuracy); (3) number of clock cycles to fetch and configure each partition (calculated based on the number of configuration words needed, which is computed with the estimation of the resources needed directly from the SUIF representation or with the number of edges, ALUs, REGs, and pre-defined values existent in the NML graph generated by MODGen).

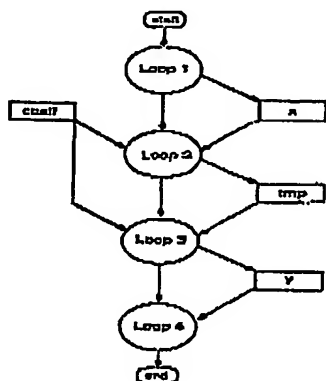


Figure 8: Top level of the HTG+ for the DCT example (this top level consists of 4 loops). Circles and boxes represent behavioral and array nodes respectively. Data is read from an input port (Loop1) and written to an output port (Loop4).

⁴E.g., twice the number of references in the xmap RAM leads to more than twice the number of objects required on XPP and delays each access because of the objects needed to MERGE and DEMUX data and address packets. Hence, combining several behavioral nodes in one partition incurs an overhead which is computed during the temporal partitioning algorithm.

The temporal partitioning algorithm starts with a partition for each node on the top of the HTG+ and then merges iteratively adjacent partitions until no performance gains are achieved considering the maximum available size for each partition. Each partition must currently define, on the control flow graph (CFG) of the program, regions of code with all entries to the same instruction and possibly multiple exists. The algorithm considers the overlapping of configuration and execution with fetch during the merging of partitions. The algorithm starts with the granularity of the nodes in the HTG+ and only if a block node cannot be mapped it considers partitioning at the statement or sub-block level. Thus, the granularity of the algorithm adapts according to the application needs.

The temporal partitioning strategy only exploits configuration boundaries inside loop bodies if an entire loop cannot be mapped to the XPP or contains more than one inner loop in the same level of the loop body. If these cases occur, the algorithm is applied hierarchically to the body of the loop.

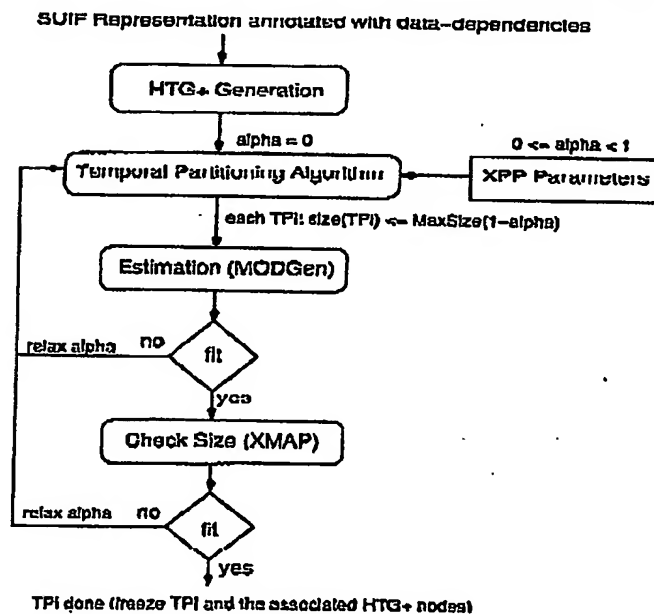


Figure 9: Automatic temporal partitioning methodology.

Fig.9 shows the methodology which uses three levels (the computational efforts increase from the first to the third level): (1) Temporal Partitioning algorithm based on the estimation of the needed resources done with function costs based on the number and kind of operations in the source code. The algorithm uses the HTG+ and the SUIF representation of the program; (2) For each configuration, selected in the first level, the estimated sizes are checked with the ones estimated by generating the NML.

graph with MODGen. If the size surpasses the available resources, the algorithm rerun level (1), relaxing the size constraint (diminishing the maximum number of available resources); (3) Check if each configuration successfully checked in level (2) can be really mapped to the XPP. This level uses functions of the mapper, placer and router. If the configuration cannot be implemented in the XPP, the algorithm returns to level (1), once more relaxing the size constraint. The size constraint is relaxed by reducing the alpha parameter in each backward iteration (see Fig.9).

After exposing the configurations with TempPart, the compiler introduces the statements needed to communicate scalar variables between partitions (see Fig.10). Arrays are used as inter-partition storage for scalar variables too, since only RAMs (to which the arrays are mapped) keep their data during reconfiguration. TempPart also ensures that arrays used by more than one configuration, or by the same configuration loaded more than once onto the XPP, are bound to the same memory location and such location is not used by other arrays during the lifetime of the array variable. The assignment of all arrays (the initially used in the source code plus the added ones to communicate data) to the internal memories is done based on the lifetimes of the arrays determined by the sequence of configurations that were previously exposed in the input program. This permits, in some cases, to use less internal memories since they can be time shared, among different configurations.

```

int comm[11];
...
a = b * c;
XPP_next_conf();
d = a/e;
...
a)

```

...	a = b * c;	#1
comm[0] = a;		#1
a = comm[0];		#2
d = a/e;		#2
b)	c)	

Figure 10: Example illustrating the communication of the value of a scalar variable between two configurations. a) source code; b) code with statements inserted to buffer the data; c) configuration ID for each of the statements in b).

4.4 Generating the NML Application

Each partition is input to MODGen, which generates the NML structure to be mapped to the XPP. MODGen generates, for each exit point existent in each partition, an event connected to one of the CM ports available in the XPP (the CM can check if an event is generated and can proceed with different configurations based on the value of the event). The compiler generates both the NML representation of each partition and the NML section specifying the control flow of configurations. Such control flow is orchestrated by the CM of the XPP during runtime, as has been

already explained.

The compiler also generates NML code considering the pre-fetch (load of a configuration to the cache of the XPP) of configurations. The compiler can furnish two different strategies: (1) request of the pre-fetch of all configurations existent in the application in the start of the execution; (2) request in each configuration of the pre-fetch of the next. The request is done before the start of the configuration step for the current configuration. Strategy (1) is used most of the times. However, there are cases where using (2) is better. In the presence of several nested if-then-else structures with different configurations for each branch, a pre-fetch sequence defined at compile time can introduce too much overhead.

5 Experimental Results

Tab 1 shows some results obtained when compiling a set of benchmarks with the XPP-VC. Note that none of the examples shown was specially coded to exploit more efficiently the architectural features of the XPP (e.g., partitioning and distribution of arrays among the internal memories) and thus the results can be further improved. An XPP Core with a single PAC was used. The 2nd column represents the size of the PAC (number of columns and rows of PAEs) used for each example. Columns #cf, #PAE, #Lat, and #max represent the number of configurations, number of PAEs used (it is shown the maximum number of PAEs of the largest configuration and the total number of PAEs virtually needed), overall latency (taken into account setup, fetching, configuration, data communication and execution), and the maximum number of objects executing per cycle respectively. The last column shows the CPU time (using a Pentium III @933MHz with Linux) to compile each example (from the source program to the generation of the binary configuration file).

DCT1 is a 8x8 discrete cosine transform implementation which is based on two matrix multiplications. The algorithm uses 6 loops for the multiplications and 2 loops to stream I/O data. It is purely sequential (no unrolling is used). Temporal partitioning improves the overall latency of DCT1 by 13% and uses 31 PAEs (without partitioning 51 PAEs are used). Thus it can use a smaller XPP core. DCT2 uses the DCT kernel of DCT1 and traverses an input image of a pre-defined size (16x16 is used). It uses 2 external memories to load/store the image and 2 internal RAMs for intermediate results and to store the coefficients. The version with 6 configurations was obtained performing temporal partitioning. Since the example has two outer loops the scheme to partitioning loops was applied (the compiler uses one configuration boundary between the two main loops of

Table 1: Results obtained with XPP-VC.

Example	XPP	fref	#PAE	#Lat (ccs)	#max	#CPU time (s)
DCT1	8x12	1	51/51	12,401	7	1.0
DCT1	6x12	4	31/79	10,957	5	0.9
DCT2	6x12	1	39/39	42,034	8	1.0
DCT2	6x12	6	41/104	40,516	8	1.2
Chen	16x16	1	64/64	16,852	22	4.8
Chen	8x16	5	56/154	10,181	24	3.3
Smooth	6x12	1	44/44	7,045	8	0.8
Haar	6x12	1	58/58	10,480	10	1.1
Haar	6x12	4	39/112	7,683	10	0.9
FIR	6x12	1	49/49	4,864	28	0.9

the DCT kernel). With this scheme, a gain of 4% in performance was achieved using 30% less PAEs. Chen is a pointer-free version, with 180 lines of C, of a DCT implementation used in JPEG. Temporal partitioning furnished an improved version: 66% in performance using 12% less PAEs. The computation and data-communication is performed in 688 clock cycles. smooth represents an image filter (16x16 image). The two inner loops (3x3 window) were annotated to be unrolled and conducted to an efficient vectorization. An overall speedup of 4 (8.6 considering only execution time) over the implementation obtained without unrolling is obtained. Additionally, 2 less PAEs are used with unrolling. Haar is an implementation of the forward 2D Haar wavelet transform. An input image of 16x16 is used. A performance gain of 36% is achieved when temporal partitioning is applied. FIR is a 1D FIR filter with 12 taps filtering 2048 samples. Even with all the overheads, 0.42 samples/cycle is computed (0.87, considering only the latency to communicate data to external memories and the FIR computation).

Each one of the examples was compiled in less than 5 seconds. This reveals that it is possible to have runtimes comparable to the ones achieved by software compilation. Performance gains obtained with temporal partitioning are shown. Since we run most examples with small data and image sequences, the configuration overhead is significant.

Note that the current methodology does not use neither the full potentialities of the XPP nor some optimizations: (1) The execution of a partition only starts after the full configuration of its resources; (2) No pipelining between fetch and configuration for the same partition has been used; (3) The capacity of the XPP to configure concurrently distinct PACs was not used; (4) An arbitrary order for pre-fetching of configurations conditionally requested is used (the order should be based on the most frequently taken path, e.g., determined by profiling); (5) The configuration FIFOs in each array-object were not used. Hence, the performance results

can be further improved.

6 Related Work

The XPP technology offers a promising reconfigurable computing platform. Being a step forward in the context of reconfigurable computing, it permits to attack some of the well-known deficiencies of related technologies. The following sub-sections illustrate the most closely related work and reveals the most important differences.

6.1 High-Level Compilation

The work on compiling high level descriptions onto reconfigurable logic has been the focus of many researchers since the first simple attempts [7]. Most of this work targets FPGA devices and thus need logic synthesis, even when module generators are included in the compilation flow, as is the case with the MARGE [8] compiler. In addition, such approaches also need backend mapping, place and route, which are very time consuming with FPGA technology. Even when pre-placed and pre routed components are used to assist the compilation flow, the compilation time is still in the order of minutes or hours.

New approaches have been used, which target research architectures. One of those approaches is the Garp-C compiler [9]. Although it is used for a reconfigware/software architecture, the configuration bit stream generation, based on exploitation of instruction-level parallelism beyond basic blocks and assisted with fast mapping and placement tasks permits to target fine-grain reconfigurable architectures efficiently with short compilation times.

As Garp-C and MARGE, XPP-VC also uses the SUIF compiler front-end. The generation of the hardware structure to be mapped to the XPP is assisted with the pipeline vectorization ideas presented in [6]. However, the generation of the control structure, based on the event packets of the XPP is completely new. Since the XPP is a coarse-grained architecture, which directly supports arithmetic and other operations occurring in high-level languages, there is no need for complex synthesis and mapping. The control structure is also directly mapped to objects handling events.

6.2 High-Level Temporal Partitioning

Temporal partitioning at the behavioral level has been already successfully conducted for FPGAs and other type of RPU. The majority of the current approaches try to use a minimum number of configurations by using all the possible RPU size available for each temporal partition (see, for instance, [10]). Such schemes only consider another partition after the current one has filled the available resources

and are insensible to the optimization that must be applied to reduce the overall execution by overlapping the fetching, configuration and execution steps. Albeit not considering such optimizations, ILP formulations presented by some authors [11] are incapable to deal with the complexity of many realistic applications.

One of the first attempts to reduce the configuration overhead in the context of temporal partitioning has been presented in [12]. However, the approach uses the simple model of splitting the available FPGA resources into two parts and performing temporal partitioning using half of the total available area as the size constraint. The scheme only overlaps configuration with execution of adjoining partitions and does not take into account the pre-fetch steps that can be efficiently used in some RPU architectures. Furthermore, the approach causes problems, when some resources of the RPU must be shared by two or more partitions. This contradicts the requirement of disjoint spaces of the RPU used by two adjacent temporal partitions.

The temporal partitioning algorithm used in the XPP-VC compiler is based on some ideas presented in [13]. The special characteristics of the algorithm to deal with resource-sharing during the creation of the partitions are not used and special heuristics have been added to deal with the fetch and configuration time of each partition. The purposed partitioning of loops was firstly introduced in this paper. The scheme can deal with any type of loops. The previous approaches consider loop distribution when a loop does not fit onto the RPU [14]. However, loops which cannot be entirely mapped onto a single configuration and which cannot be distributed are not compiled. Our method can deal with programs with unlimited complexity as long as the supported C subset is used. It does not depend on the feasibility of a specific compiler transformation.

7 Conclusions and Future Work

This paper describes the new Vectorizing C Compiler, XPP-VC, which maps programs in a C-subset extended by port access functions to PACT's XPP architecture. Assisted with a fast place and route tool, it furnishes a complete "push-button" path from algorithmic descriptions onto XPP configuration data with short compilation times.

An innovative temporal partitioning scheme is presented. It enables the mapping of complex programs and furnishes XPP applications with performance gains by hiding some of the configuration time. A new mechanism to handle partitioning of loops, which supports loop execution by the configuration manager of the XPP, is also presented. Furthermore, the compiler generates self-contained configuration data even when several configurations are exposed.

Ongoing work focuses on tuning the estimation steps to assist automatic temporal partitioning and on improving the configuration data generated.

In addition to loop unrolling, *loop merging*, *loop distribution* and *loop tiling* will be used to improve loop handling, i.e., enable more parallelism or better XPP usage. A future extension of the compiler for a host-XPP hybrid system is planned. The compiler will map suitable program parts, especially inner loops, to the XPP, and the rest of the program to the host processor.

8 Acknowledgements

We gratefully acknowledge the support of the PACT team, especially the support and assistance of Armin Strobl, Daniel Bretz, and Frank May related to the XDS tools.

References

- [1] R. Hartenstein, "A Decade of Reconfigurable Computing: a Visionary Retrospective," In *Proc. Design, Automation and Test in Europe (DATE'01)*, 2001, pp. 642-649.
- [2] PACT Informationstechnologie GmbH, Germany, "The XPP White Paper," Release 2.0, June 2001, <http://www.pactcorp.com>
- [3] V. Baumgarte, F. May, A. Nüchel, M. Vorbach, M. Weinhardt, "PACT XPP - A Self-reconfigurable Data Processing Architecture," In *Proc. Int'l Conf. on Engineering of Reconfigurable Systems and Algorithms (ERSA'01)*, Las Vegas, Nevada, USA, June 25-28, 2001, pp. 64-70.
- [4] SUIF Compiler system, "The Stanford SUIF Compiler Group," <http://suif.stanford.edu>
- [5] Muchnick, S. S., *Advanced Compiler Design and Implementation*, Morgan Kaufmann Publishers, Inc., San Francisco, CA, USA, 1997.
- [6] M. Weinhardt and W. Luk, "Pipeline Vectorization," In *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, Feb. 2001, pp. 234-248.
- [7] Ian Page and Wayne Luk, "Compiling occam into FPGAs," In *FPGAs*, Will Moore and Wayne Luk, editors, Abingdon FF&CS Books, Abingdon, England, UK, 1991, pp. 271-283.
- [8] Maya Gokhale, and Edson Gomersall, High-Level Compilation for Fine Grained FPGAs, In *Proc. IEEE 5th Symposium on Field-Programmable Custom Computing Machines (FCCM'97)*, Napa Valley, CA, USA, April 16-18, 1997, pp. 165-173.
- [9] T. J. Callahan, J. R. Hauser, and J. Wawrzynck "The Garp architecture and C compiler," In *IEEE Computer*, 33(4), April 2000, pp. 62-69.

- [10] João M. P. Cardoso, and Horácio C. Neto, "An Enhanced Static-List Scheduling Algorithm for Temporal Partitioning onto RPLs," In *Proc. IFIP X Int'l Conf. on Very Large Scale Integration (VLSI'99)*, Lisbon, December 1-3, 1999. Kluwer Academic Publishers, pp. 485-496.
- [11] I. Ouass, et al., "An Integrated Partitioning and Synthesis System for Dynamically Reconfigurable Multi-FPGA Architectures," in *Proc. 5th Reconfigurable Architectures Workshop (RAW'98)*, Orlando, Florida, USA, March 30, 1998, pp. 31-36.
- [12] Satish Ganesan, and Ranga Vemuri, "An Integrated Temporal Partitioning and Partial Reconfiguration Technique for Design Latency Improvement," in *Proc. Design, Automation & Test in Europe (DATE'00)*, Paris, France, March 27-30, 2000, pp. 320-325.
- [13] João M. P. Cardoso, "A Novel Algorithm Combining Temporal Partitioning and Sharing of Functional Units," In *Proc. IEEE 9th Symposium on Field-Programmable Custom Computing Machines (FCCM'01)*, Rohnert Park, CA, USA, April 30 - May 2, 2001.
- [14] Meenakshi Kaul, Ranga Vemuri, Sriram Govindarajan, Iyad E. Ouass, "An Automated Temporal Partitioning and Loop Fission approach for FPGA based reconfigurable synthesis of DSP applications," in *Proc. IEEE/ACM Design Automation Conference (DAC'99)*, New Orleans, LA, USA, June 21-25, 1999, pp. 616-622.

Temporal Partitioning for the XPP-VC Compiler

<http://www.pactcorp.com>

18th January 2002

1 Benefits of Temporal Partitioning

Temporal partitioning can have a distinct and important goal than to simply enable the compilation of algorithms which the mapping onto the RPU (Re-configurable Processing Unit) resources cannot be accomplished by only one configuration. For instance, temporal partitioning targeting the XPP [1] can reduce, when efficiently applied, the overall execution latency. Such reduction can be mainly enabled by the following issues:

- reduction of the interconnection lengths, by reducing each design complexity, can furnish better performance results (long interconnections pass through registers and thus adding clock cycle delays);
- reduction of each temporal partition complexity can itself reduce the number of registers used for vertical routing;
- reduction of the number of references, in each temporal partition, using the same resource, by distributing the overall references among temporal partitions, can furnish better performance results as well. This happens with the statements presented in the program referring the same array;
- reduction of the overall configuration overhead by overlapping fetching, configuration and execution of distinct temporal partitions.

The reduction of the configuration overhead is due to 3 distinct sources of overlapping, possible with the XPP architecture:

1. loading of subsequent configurations into the cache in parallel with the configuration of the current one;
2. execution of one configuration while the next one is being configured;
3. execution of one configuration while the next one is being loaded into the cache.

2 SPECIFICATION OF CONFIGURATIONS

(2) can have more impact if the map, place, and route phase could, as most as possible, to confine temporally adjacent configurations in distinct locations of the XPP (augmenting the concurrency between execution and configuration).

For reconfigurable computing platforms, where the reconfiguration of the array takes several clock cycles, is in most of the cases preferable that a configuration is reused as much time as possible in order to reduce the reconfiguration overhead. Thus, loops in the source code are always good candidates to be entirely implemented in a single configuration.

2 Specification of Configurations

Configurations can be specified by the programmer using *XPP_next_conf()* statements in the source code of a given application. Such statements must expose, on the control flow graph (CFG) of the procedure, regions of code with all entries to the same instruction and eventually multiple exists. The compiler exposes the configurations, removes such statements from the SUIF1[2] intermediate representation, checks for invalid specifications of configuration boundaries (when the statements expose regions with entries to different statements in a region of code, or when code can be contained in more than one region¹), inserts the code responsible to the data communication between temporal partitions, and generates both the NML (Native Mapping Language) [8] representation of each configuration and the application section specifying the control flow of configurations. Such control flow is orchestrated by the CM (Configuration Manager) of the XPP during runtime.

Consider a pointer-free version of the quantizer of an h263 implementation [3] which code is shown in Fig. 1. Four *XPP_next_conf()* statements were inserted in the code to specify three configurations. The configurations specified are represented in the CFG of the example that can be seen in Fig. 2. Apart from specifying temporal partitions in such a way that the mapping to XPP is accomplished, there can be the case that, merging only the mostly taken conditional paths in the same configuration can reduce the total execution time by substantially reducing the reconfiguration time (since the partitions for the other paths are not configured when they are not taken). Fig. 2 presents such a case. If the path bb_0, bb_1 and bb_2 was identified as the most frequently executed, such path can be specified to be in the same configuration². In such a case, the configurations related to bb_3 and bb_4 will only be called when the most frequently path has not been taken. In some examples, paths are only executed in "debug mode" (as is the case of the branch taken when QP evaluates to false in the source code of Fig. 1).

After exposing the configurations, the temporal partitioning phase introduces the statements needed to communicate scalar variables between two different configurations (see Fig. 3). Currently, the scalar variables are stored in

¹Tail duplication could be applied in some examples.

²Tail duplication of bb_5 would permit to have a configuration with {bb_0, bb_1, bb_2, bb_3}; another one with {bb_4, bb_5}; and another one with {bb_3, bb_5}.

João M P Cardoso PACT Informationstechnologie GmbH,
http://www.pactcorp.com, December 18th, 2001.

2 SPECIFICATION OF CONFIGURATIONS

```

if (QP) {
    if (Mode == MODE_INTRA || Mode == MODE_INTRA_Q) { /* Intra */
        qcoeff[0] = mmax(1, mmin(254, coeff[0]/8));
        for (i = 1; i < M; i++) {
            level = (abs(coeff[i])) / (2*QP);
            qcoeff[i] = mmin(127, mmax(-127, sign(coeff[i]) * level));
        }
        XPP_next_conf();
    } else { /* non Intra */
        XPP_next_conf();
        for (i = 0; i < M; i++) {
            level = abs(coeff[i]) - QP/2 / (2*QP);
            qcoeff[i] = mmin(127, mmax(-127, sign(coeff[i]) * level));
        }
        XPP_next_conf();
    }
} else {
    XPP_next_conf();
    /* No quantizing.*/
    for (i = 0; i < M; i++) {
        qcoeff[i] = coeff[i];
    }
    XPP_next_conf();
}
}

```

Figure 1: C source code of the quantization algorithm with configuration boundaries specified.

2 SPECIFICATION OF CONFIGURATIONS

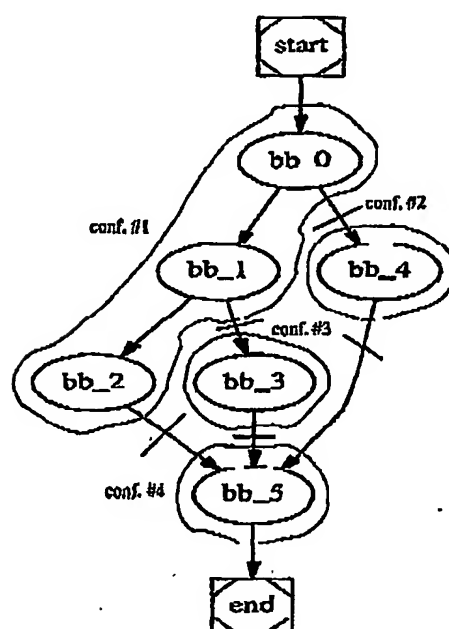


Figure 2: CFG of the algorithm shown in Fig. 1. The lines crossing edges represent the *XPP_next_conf()* statements in the code. The bubbles containing basic blocks of the CFG represent the exposed regions of the CFG that are implemented in different temporal partitions.

João M P Cardoso/PACT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

3 MAPPING APPLICATIONS WITH CONFIGURATION BOUNDARIES IN LOOP BODIES

	int comm[1];	
...	...	conf #1
a = b * c;	a = b * c;	conf #1
next_conf();	comm[0] = a;	conf #1
d = a/e;	a = comm[0];	conf #2
...	d = a/e;	conf #2
a)	b)	c)

Figure 3: Example illustrating the communication of the value of a scalar variable between two configurations. a) source code; b) source code with statements inserted to buffer the data; c) configuration ID for each of the statements in b).

arrays specially inserted in the SUIF1 representation of the given application. Those arrays are mapped to internal memories of the XPP³. The temporal partitioning phase also ensures that arrays used by more than one configuration, or by the same configuration loaded more than once to the XPP, are binded to the same memory location and such location is not used by other arrays during the lifetime of the array variable⁴.

The assignment of the overall existent arrays (the initially used in the source code more the added ones to communicate data) to the internal memories is done based on the lifetimes of the arrays determined by the sequence of configurations that were previously exposed in the application. This permits, in some cases, to reduce the number of internal memories needed by time sharing, among different configurations, some internal memories during the execution of the application on the XPP.

The XPP-VC compiler generates, for each exit point existent in each configuration, an event connected to one of the CM ports available in the XPP (the CM can check if an event is generated and can proceed with different configurations based on the value of the event). The generated event has value "0" if the path that activates that exit is taken and "1" otherwise.

3 Mapping Applications with Configuration Boundaries in Loop Bodies

Configuration boundaries in loop bodies can be deal by performing loop distribution (as long as it can be applied) or by temporal partitioning the loop and

³In this case, the internal memories are used as data buffers for the maintenance of the original program behavior.

⁴At the moment each configuration must use a number of array variables, to be assigned to the internal memories of the XPP, less or equal than the number of internal memories of the XPP (the compiler assigns each array to a distinct memory). However, the total number of arrays existent on the overall configurations can surpass the number of internal memories, if some memories can be shared between configurations due to the non-overlap of the lifetime of array variables. The data stored in memories is maintained across reconfigurations.

3 MAPPING APPLICATIONS WITH CONFIGURATION BOUNDARIES IN LOOP BODIES

int i;	int i;	
...	i=0;	conf #1
for(i=0;i<N;i++) {	lab1: if(i<N) {	conf #2
statement1;	statement1;	conf #2
XPP_next_conf();	statement2;	conf #3
statement2;	i++;	conf #3
}	goto lab1;	conf #3
statement3;	}	conf #3
	statement3;	conf #4
a)	b)	c)

Figure 4: Example of the transformation applied to loops with configuration boundaries in their bodies. a) original source code; b) transformed code; c) configuration ID for each statement in b).

use the CM to orchestrate the control flow.

Currently, loop distribution [11] is not automatically applied⁵. All the loops with configuration boundaries specified in their bodies are transformed into *if()* *goto label* kind loops in order to permit the NML generation by the XPP-VC compiler. Fig. 4 shows an example of such transformation without the statements needed to communicate the value of scalar variables between configurations. The 3rd column shows the configuration ID of each statement. Each configuration requests the next configuration to be taken (if the exit taken is to the end then only the "reconf"⁶ of the configuration is done). Conf #2 needs a conditional request mechanism to call conf #3 or conf #4 based on the value of the *i<N* expression. Since conf #3 always requests, at the end of its execution, conf. #2, the initial behavior of the loop is maintained. The temporal partitioning task also creates two more configuration boundaries to preserve the initial functionality. From Fig. 4b) can be seen that configuration boundaries were inserted before and after the *if* statement. Such boundaries are needed since the code before and after will be executed once and both the *if* header and body will iterate *N+1* and *N* times respectively.

The configuration boundaries inserted in loop bodies must specify, at the scope of the loop body, the permitted type of regions (already explained).

Loop distribution (also known as "loop fission") will be the preferable form to implement loops, which generated NML does not entirely fit in the available resources of the XPP. Such transformation can potentially lead to the introduction of temporary arrays. Consider the loop shown in Fig. 5 where a configuration boundary is specified. The loop can be splitted so that the two statements are each one in one loop and the configuration boundary is now outside any loop body. However, we need to scalar expand variable *s* in order to maintain the

⁵The compiler should check if the loop distribution can be applied on each temporal partition boundary existent in loop bodies.

⁶A reconf means that the resources used by that configuration are released and then can be reconfigured.

4 EXECUTION STRATEGIES

```

...
for(i=0; i<N; i++) {
    s = ... // statement 1
    XPP_next_conf();
    ... = s // statement 2
}
...

int tmp_s[N];
...
for(i=0; i<N; i++) {
    tmp_s(i) = ... // statement 1
}
XPP_next_conf();
for(i=0; i<N; i++) {
    ... = tmp_s(i) // statement 2
}

```

Figure 5: Applying loop distribution as another way to enable temporal partitioning on loop bodies.

initial functionality (one array with the size of the number of iterations of the loop must be declared and is used to communicate each *s* value in each of the loop iterations).

4 Execution Strategies

To reduce the overall latency, efficient exploitation of the pipelining of the 3 steps presently in each temporal partition (fetching, configuring, and array execution) must be conducted.

Two "reconf" modes can be used (the user can select one of the modes in the options of the XPP-VC compiler related to temporal partitioning):

- "reconf" executed by the CM. In this case each configuration communicates with the CM sending an event, when the completion of execution, to request the next configuration. This next configuration starts by executing a "reconf" command to an XPP resource of the configuration (that command is broadcasted throughout all the resources used by the configuration, and so the resources will be released and can be reconfigured by the next configuration). When a configuration can be requested by more than one previous configuration, special configurations are inserted in the Temporal Partition Control Flow Graph (TPCFG⁷) between each source and the sink. Such special configurations only command the "reconf" of a resource in the XPP of the previous configuration and request the next one. This type of "reconf" does not permit to have overlapping between execution and configuration between temporal partitions;
- "reconf" self applied by each configuration. In this case each configuration at the end of the execution broadcasts a "reconf" event to all the XPP

⁷The TPCFG is a directed, eventually cyclic, graph where each node represents a configuration (temporal partition) and each edge between two nodes specifies the execution flow of the application through its temporal partitions. There is only one edge between two nodes of the graph and each node represents a region of the CFG of the application.

João M P Cardoso, PACT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

5 AUTOMATIC TEMPORAL PARTITIONING

resources pertaining to it. This mode does not need addition of special configurations by the compiler and permits that the CM try to configure the next configuration called during the execution of the callee temporal partition (when only one configuration path is presented).

The compiler also generates NML code considering the pre-fetch (load of a configuration to the cache of the XPP) of configurations. When the pre-fetch is enabled two strategies can also be automatically used:

- request of the pre-fetch of all configurations existent in the application in the start of the execution (during this pre-fetching the flow of configuration and execution is done in the way it is specified in the application section of the NML file);
- request in each configuration of the pre-fetch of the next. The request is done before the start of the configuration step for the current configuration.

The CM of the XPP permits also speculative configuration of a temporal partition that can conduct to better performance results even when the Map, Place and Route does not try to locate temporal partitions in non-overlapping areas of the XPP. The strategy tries to configure the partition speculatively used after the configuration of the current one. If the path which includes that configuration is taken, the CM only has to enable the start of the execution of the configuration (see the section of the NML code in Fig. 6 and the simulation results in Fig. 7, where conf_MOD2 is speculatively configured during the execution of conf_MOD0). When such path is not taken, the CM releases the resources already configured and requests the other configuration.

5 Automatic Temporal Partitioning

Automatic temporal partitioning permits the automatic exposing of configurations oriented by two distinct goals:

- minimum number of configurations: this goal can be achieved with algorithms that try to use all the available reconfigurable processing units during the assignment of segments of behavioral code to the same configuration;
- minimum overall latency: this goal can be achieved by considering the costs to load into the cache, to configure and to execute each configuration with the XPP array. An important strategy that must be considered is the use of pre-fetch of configurations while one of the others is running. Arrays of constants or with pre-defined values used in one or more configurations can be initialized in one of the previous configurations if such one exists. This takes advantage of the initialization of the array carried out by using the configuration bus.

João M P Cardoso SPAUT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

5 AUTOMATIC TEMPORAL PARTITIONING

```

...
CONFIG conf_MOD0 {
  CONF_MODULE(MOD0)           // request the configuration of MOD0
  SET(MOD0.Reconf.E = 1)      // enable the self releasing of resources
  REQUEST(conf_MOD2_spec)     // speculative configuration
  // if (MOD0.CMPort0 == "0") then conf_MOD2_exec is requested
  // else continue
  // if (MOD0.CMPort1 == "Q") then conf_MOD1 is requested
  CONF_CMPORT(MOD0.CMPort0, conf_MOD2_exec, _) // take MOD2 or continue
  CONF_CMPORT(MOD0.CMPort1, conf_MOD1, _)      // take MOD1
}
CONFIG conf_MOD2_spec {
  CONF_MODULE(MOD2)           // request the configuration of MOD2
}
CONFIG conf_MOD2_exec {      // MOD2 is taken
  SET(MOD2.Start.A = 1)      // enable the start of computing of MOD2
  SET(MOD2.Reconf.E = 1)     // enable the self release of resources
  REQUEST(conf_MOD3)         // request the next configuration
}
CONFIG conf_MOD1 {           // MOD1 is taken
  REQUEST(conf_MOD2_rec)     // request the releasing of resources
  CONF_MODULE(MOD1)          // request the MOD1
  SET(MOD1.Reconf.E = 1)     // enable the self releasing of resources
  REQUEST(conf_MOD3)         // request the next configuration
}
CONFIG conf_MOD2_rec {
  RECONF(MOD1.Start)         // release the resources of MOD1
}
...

```

Figure 6: Example of a section of NML code describing the speculative configuration concept.

5 AUTOMATIC TEMPORAL PARTITIONING

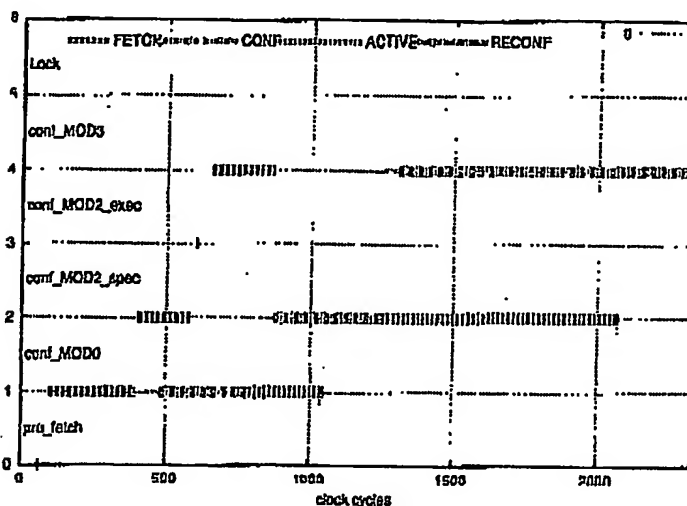


Figure 7: Example of the overlapping among the fetch, configuration, and execution steps of different temporal partitions.

From the SUIF1 representation of the C source code the temporal partitioning phase constructs an extended HTG (Hierarchical Task Graph)⁸. Such extended graph has 2 types of nodes:

1. behavioral nodes representing source lines of code in the input program;
2. array nodes representing each array existent in the source code.

Type (1) nodes have 3 distinct sub-types:

1. block nodes representing basic blocks with one-entry and a single exit;
2. compound nodes representing if-then-else structures;
3. loop nodes representing the loops (for, while, ect.). Loop and compound nodes explicitly embody hierarchical levels.

Edges in the HTG+ represent data communication between two nodes or just enforce execution's precedence.

Each behavioral node of the HTG+ is labeled with the following information (some of the labelling steps require estimation efforts):

- block and compound nodes: number of ALUs and REGs;

⁸The model has been chosen, because it will also permit to exploit loop and task level parallelism.

João M P Cardoso/PACT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

5 AUTOMATIC TEMPORAL PARTITIONING

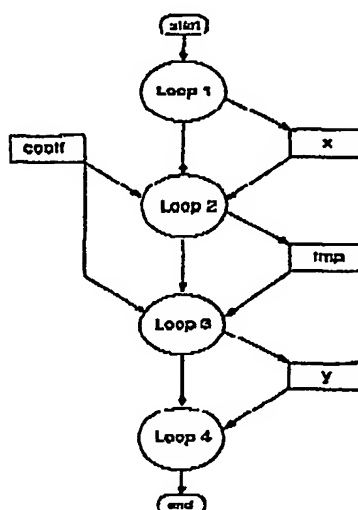


Figure 8: Top level of the HTG+ for the DCT example (this top level consists of 4 loops). Circles and boxes represent behavioral and array nodes respectively. Loop 1 reads the data-stream from an input port to an internal memory and Loop 4 writes the data-stream generated by the DCT code (Loops 2 and 3) from an internal memory to an output port of the XPP.

- loop nodes: number of iterations (unknown if unbound), and number of ALUs and REGs;
- array nodes: the size of the array, type of the elements, and, when they do exist, the initialization values.

Each edge between two behavioral nodes of the HTG+ is labeled with the number of data words that must be transferred between the two nodes.

Each edge between an array and a behavioral node in the HTG+ is labeled with the number of load and store references ($\dots = A[i]$ and $A[i] = \dots$ respectively) in the source code represented by the behavioral node to that particular array. The estimated number of times that each load and store reference will be executed is also collected. Such information is used to calculate the penalty when two or more behavioral nodes are merged into the same temporal partition. Such penalty is related to the use of the same array by different behavioral nodes and adds an overhead to the execution latency of that temporal partition and to the number of resources needed for its implementation.

Fig. 8 shows the top level of the HTG+ for an implementation of the DCT (Discrete Cosine Transform) based on matrix multiplications.

The automatic temporal partitioning phase needs 3 types of estimations:

João M P Cardoso | PACT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

5 AUTOMATIC TEMPORAL PARTITIONING

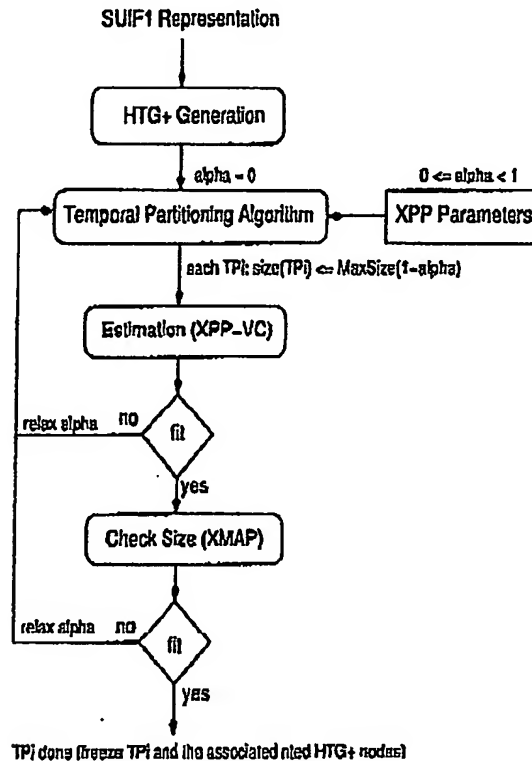


Figure 9: Automatic Temporal Partitioning methodology.

- number of XPP resource units needed by the configuration implementing a single or a set of behavior nodes;
- latency for a behavior node or a set of connected behavior nodes on the HTG+ (this does not need to be accurate to the real execution time and only needs to have relativeness accuracy);
- number of clock cycles to fetch and configure each temporal partition (calculated based on the number of configuration words needed⁹).

The temporal partitioning strategy does not exploit configuration boundaries inside loop bodies, unless the entire loop cannot be mapped to the XPP. The generation of this type of temporal partitions never produces better results

⁹Can be estimated by the number of edges, ALU nodes, REG nodes, and pre-defined values existent in the hardware graph generated by the XPP-VC compiler.

João M P Cardoso, PACT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

6 DISCUSSION

(at least when the loop behavior is ensured by the CM). The justification is supported by the reutilization of resources, already configured, achieved when the entire loop is implemented by a single configuration. When a loop does not fit in the XPP, the algorithm is applied hierarchically to the body of the loop.

Fig. 9 shows the methodology. The strategy works around 3 levels (the computational efforts increase from the first to the third level):

1. Temporal Partitioning algorithm based on the estimation of the needed resources done with function costs based on the number and kind of operations in the source code. The algorithm uses the HTG+ and the SUIF representation of the program;
2. For each configuration, selected in the first level, the estimated sizes are checked with the ones estimated by generating the NML graph with the XPP-VC compiler. If the size surpasses the available resources, the algorithm rerun level one, relaxing the size constraint (diminishing the maximum number of available resources);
3. Check if each configuration successfully checked in level 2 can be really mapped to the XPP. This level uses functions of the mapper, placer and router. If the configuration cannot be implemented in the XPP, the algorithm returns to level 1, once more relaxing the size constraint.

The temporal partitioning algorithm used is based on the ideas presented in [7]. The special characteristics of the algorithm to deal with resource-sharing during the creation of the temporal partitions have been removed and special heuristics have been added to deal with the fetch and configuration time of each temporal partition. The algorithm tries to overlap configuration and execution with fetch during the selection of the HTG+ nodes to each temporal partition.
[describe more]

6 Discussion

We call the attention of the reader for the fact that the current methodology does not use neither the full potentialities of the XPP nor some optimizations:

1. The execution of a given temporal partition only starts after all the used resources have been configured;
2. No pipelining between fetch and configuration has been used. The configuration of the XPP resources for a specific temporal partition only starts after its configuration words are fetched (loaded to the XPP cache);
3. No overlapping on execution between two or more configurations;
4. The capacity of the XPP technology to configure concurrently distinct PACs (each PAC has its own CM);

João M P Cardoso, PACT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

7 RELATED WORK

5. An arbitrary order for fetching of temporal partitions conditionally requested is used (the fetching order should be done by the most taken path determined by profiling);
6. Behavioral nodes exposed in the HTG+ as concurrent nodes are not at the moment implemented, by the XPP-VC compiler, with parallel execution.

Thus, we strong believe that there still be potential to improve the performance results achieved when using XPP-VC.

7 Related Work

The XPP technology offers an unique reconfigurable computing platform supported by tools that permit to compile algorithms in C. Being a step forward in the context of the reconfigurable computing it permits to attack some of the well deficiencies presently in many, if not all, other reconfigurable computing technologies. However, some of the work being done to augment the potential of such technology has sources in some works previously done.

Temporal partitioning has been already successfully conducted for FPGAs and other type of RPUs. The majority of the current approaches try to use a minimum number of configurations by using all the possible RPU size available for each temporal partition (see, for instance, [4]). Such schemes only consider another temporal partition after the current one has fulfilled the available resources and are insensible to the optimization that must be applied to reduce the overall execution by overlapping the fetching, configuration and execution steps. Albeit not considering such optimizations, ILP formulations presented by some authors [5] are incapable to deal with the complexity of many realistic examples.

One of the first attempts to reduce the configuration overhead in the context of temporal partitioning has been presented in [6]. However, the approach uses the simple model of splitting the available FPGA resources into two parts and performing temporal partitioning using half of the total available area as the size constraint. The scheme only overlaps configuration with execution of adjoining partitions and does not enter into account to the pre-fetch steps that can be efficiently used in some RPU architectures. Furthermore, the approach can originate some problems, when some resources of the RPU must be shared by two or more partitions (eliminating the requirement of disjoint spaces of the RPU used by two adjacent temporal partitions).

[12] presents the scheduling of kernels (sub-tasks) targeting the Morphosys architecture. They use an efficient search pruning scheme added to an heuristic that permits to consider firstly solutions which potentially conduct to the best performance results. However, they mainly orient the search to data re-use among the scheduled kernels which is only suitable to type of reconfigurable computing architectures where no local memories to the RPU are available. The scheduler tries to overlap computing and data transfers and minimize context reloading, which as we can see from the examples shown can not always conduct

João M P Cardoso / FACT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

REFERENCESREFERENCES

to the overall minimum latency. The scheme needs as input the application flow graph (without concurrency and conditional paths) and the kernel timing. The approach does not consider temporal partitioning and so needs that each kernel configuration does not exceed the context memory size.

References

- [1] V. Baumgarte, F. May, A. Nuckel, M. Vorbach, M. Weinhardt, "PACT XPP A Self reconfigurable Data Processing Architecture," In *Proc. of the International Conference on Engineering of Reconfigurable Systems and Algorithms (ERSA'01)*, Monte Carlo Resort, Las Vegas, Nevada, USA, June 25-28, 2001, pp. 64-70.
- [2] SUIF Compiler system, "The Stanford SUIF Compiler Group," <http://suif.stanford.edu/>
- [3] <http://kimp.sarang.net/text/h263/>
- [4] João M. P. Cardoso, and Horácio C. Neto, "An Enhanced Static-List Scheduling Algorithm for Temporal Partitioning onto RPU's," In *Proc. of the IFIP TC10 WG10.5 X International Conference on Very Large Scale Integration (VLSI'99)*, Lisbon, December 1-3, 1999. VLSI: Systems on a Chip, Luis M. Silveira, Srinivas Devadas and Ricardo Reis (Editors), Kluwer Academic Publishers, pp. 485-490.
- [5] I. Onaiss, et al., "An Integrated Partitioning and Synthesis System for Dynamically Reconfigurable Multi-FPGA Architectures," in *Proc. of the 5th Reconfigurable Architectures Workshop (RAW'98)*, Orlando, Florida, USA, March 30, 1998, pp. 31-36.
- [6] Satish Ganesan, and Ranga Vemuri, "An Integrated Temporal Partitioning and Partial Reconfiguration Technique for Design Latency Improvement," in *Proc. of Design, Automation & Test in Europe (DATE'00)*, Paris, France, March 27-30, 2000, pp. 320-325.
- [7] João M. P. Cardoso, "A Novel Algorithm Combining Temporal Partitioning and Sharing of Functional Units," In *Proc. of the IEEE 9th Symposium on Field-Programmable Custom Computing Machines (FCCM'01)*, Rohnert Park, California, USA, April 30 - May 2, 2001, IEEE Computer Society Press, Los Alamitos, CA, USA.
- [8] PACT Informationstechnologie GmbH, Germany, "NML Reference: Release 2.0," April 2001.
- [9] PACT Informationstechnologie GmbH, Germany, "KDS ...".
- [10] PACT Informationstechnologie GmbH, Germany, "XPP-VC...".

João M P Cardoso, PACT Informationstechnologie GmbH,
<http://www.pactcorp.com>, December 18th, 2001

REFERENCESREFERENCES

- [11] Muchnick, S. S., *Advanced Compiler Design and Implementation*, Morgan Kaufmann Publishers, Inc., San Francisco, CA, USA, 1997.
- [12] Rafael Maestre, F Kurdahi, N Dagherzadeh, H. Singh, R. Hermida, and N. Fernandez, "Kernel Scheduling in Reconfigurable Computing," In *Proceedings of Design and Test in Europe (DATE'99)*, Munich, Germany, March 1999, pp. 90-96.

Abstract

Resource virtualization on FPGA devices, achievable due to its dynamic reconfiguration capabilities, provides an attractive solution to save silicon area. Architectural synthesis for dynamically reconfigurable FPGA-based digital systems needs to consider the issue of reducing the number of temporal partitions (reconfigurations), by enabling sharing of some functional units in the same temporal partition. This paper proposes a novel algorithm for automated datapath design, from behavioral input descriptions (represented by a dataflow graph), which simultaneously performs temporal partitioning and sharing of functional units. The proposed algorithm attempts to minimize both the number of temporal partitions and the execution latency of the generated solution. Temporal partitioning, resource sharing, scheduling, and a simple form of allocation and binding are all integrated in a single task. The algorithm is based on heuristics and on a new concept of construction by gradually enlarging timing slots. Results show the efficiency and effectiveness of the algorithm when compared to existent approaches.

1 Introduction

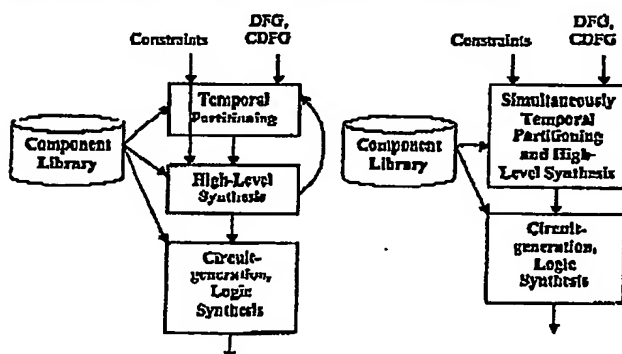
The availability of multi-programmable logic devices (such is the case of FPGAs - field programmable gate arrays) with lower reconfiguration times has made possible the concept of "virtual hardware" [1][2]: the hardware resources are supposed unlimited and implementations that oversize the resources available on the device are resolved by temporal partitioning. Then, the temporal partitioned solution is executed by time-sharing the device such that the initial functionality is preserved. This concept promises to be an efficient solution to save silicon area [1]. One of the applications is the switch among functionalities that have mutual exclusiveness on the temporal domain, such as the context-switching between coding/decoding schemes in communication, video or audio systems.

Although, even the latest commercial FPGAs, such as the XilinxTM Virtex family [3], do not have mechanisms to implement efficiently temporal partitioned functionalities and the time of reconfiguration of the overall FPGA is still quite high, the importance of the "virtual hardware" concept has already been demonstrated with computationally complex applications [4]. Industrial efforts are under way to further improve the capability of the devices to handle multiple-configurations by storing several on-chip configurations and permitting the switch between contexts in few nanoseconds [5].

The virtualization of FPGA resources has been considered by several authors while dealing with circuit netlists that oversize the available resources on the device ([6][7], just to name a few). From the point of view of the design, those approaches work at a much low-level of abstraction, without the possibility to exploit tradeoffs between the number of reconfigurations and the resource sharing of functional units (FUs), for instance. The design automation for FPGA-based systems should include temporal partitioning algorithms able to efficiently exploit the new concept. Tradeoffs among parallelism, communication costs, execution and reconfiguration times, and sharing of some FUs in the same reconfiguration need to be considered during the architectural synthesis phases.

Sharing of FUs among operations is a technique to reuse a single configuration of an FU by more than one operation of the same type. On the other hand, temporal partitioning is a technique tailored to reuse the available resources by different circuits (configurations) with the time-multiplex of the device. The nodes of a given intermediate representation (e.g., a dataflow graph) representing operations have to be scheduled in time steps to be executed in each temporal partition (TP). Temporal partitioning must preserve the dependencies among nodes (that are already temporal dependencies) such that a node B dependent on node A cannot be mapped to a partition executed before the partition where node A is mapped. In addition, considering sharing FUs during temporal partitioning can conduct to better overall results (lower number of TPs and better performance).

Figure 1a) shows a design flow which integrates temporal partitioning prior to the high-level synthesis tasks [8]. The majority, if not all, of the existent approaches utilizes the presented flow [9][10]. Our efforts address architectural synthesis¹ integrating temporal partitioning and this paper presents a new temporal partitioning algorithm that effectively takes into account sharing of FUs, while maintaining a small computational complexity. Besides, it is sufficiently flexible to target different FPGA devices. Figure 1b) shows the design flow proposed in this paper, where temporal partitioning is integrated in the high-level synthesis tasks and is performed simultaneously.



a)

b)

Figure 1. Design flow based on high-level synthesis for reconfigurable computing systems: a) traditional flow; b) proposed flow.

Example 1. Motivational example.

Consider the dataflow graph exhibited in Figure 2 (Ex1). It consists of 4 additions and 2 multiplications. Suppose that each adder uses 1 cell and has a latency of 1 clock cycle, each multiplier uses 2 cells and has a latency of 2 clock cycles and the maximum resources available on the device equals 3 cells. The dataflow graph has a critical path latency of 4 cycles and needs 8 cells given those FUs (last row of Table 1). Figure 2 shows an optimal solution (not considering the area of multiplexers, registers and control unit needed to implement sharing of a specific FU) for the example with results shown in the second row of Table 1. In Figure 2 each gray region identifies operations that are mapped to the same FU. The optimal solution is achieved with only one adder and one multiplier and fits totally on a single TP. When not considering sharing of adders, the optimum result is shown in the third row of Table 1. The algorithm proposed in this paper achieves those optimal re-

¹ There is no distinction among the terms: *high-level synthesis*, *architectural synthesis* and *behavioral synthesis*.

sults. The fourth row of the table shows the solution obtained when considering a leveling temporal partitioning algorithm that does not consider resource sharing of FUs. From this example, it can be seen that resource sharing can reduce the number of reconfigurations and can also reduce the overall execution latency. There are also cases where the critical path latency of the input dataflow graph (last row) is maintained (second row).

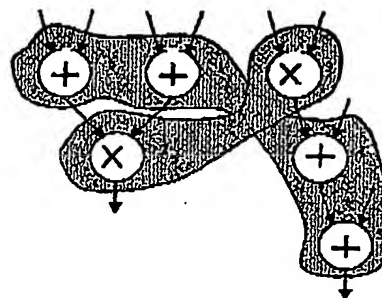


Figure 2. Dataflow graph of the example Ex1.

Table 1. Results for Ex1.

Approach	(+, *)	#TPs	Execution latency	Resources used
Optimum (sharing of adders and multipliers)	-	1	4	3
Optimum (sharing of multipliers)	-	3	5	3
ASAP (no sharing) [11]	(4, 2)	4	6	3
Without Temporal Partitioning (no sharing)	(4, 2)	-	4	8

The remainder of this paper is organized as follows. Section 2 formulates and explains the problem. The algorithm is deeply explained in section 3, where the pseudo-code and the overall performed steps are fully elucidated through an example. In section 4 experimental results are shown and discussed. In section 5, related work is described. Finally, in section 6, conclusions are presented and further work is envisaged.

2 Problem Definition

Given a dataflow graph (DFG), representing a behavioral description, $G = (V, E)$, topologically ordered, directed and acyclic, with $|V|$ nodes, $\{v_1, v_2, \dots, v_{|V|}\}$ and $|E|$ edges, where each node v_i represents an operation and each edge $e_{ij} \in E$ represents a dependence between nodes v_i and v_j . A dependence can be a simple precedence dependence or a transport-dependence due to the transport

of data between two nodes. The DFG can be obtained from an algorithmic input description. Such pre-processing step is beyond the scope of this article, but the front end of our Java compiler for reconfigurable computing systems can be employed [12].

Here we assume that there is a component library with a set of FUs and there is one FU for each type of operation in the DFG. Φ represents the set of FUs, from the component library, to be instantiated by the algorithm. R_{MAX} represents the resource capacity available on the device, $R(\pi_i)$ returns the number of resources utilized by the TP π_i and $R(v_i)$ returns the number of resources utilized by the FU instance associated with v_i . $N(\pi_i)$ returns a subset of nodes of V mapped to π_i .

Each partition π_i is a non-empty subset of V , where for each node exists a map to one and only one FU instance in Φ . $\pi(v_i)$ identifies the TP where node v_i is mapped. The set of the TPs is represented by:

$$\wp = \bigcup_{i=1}^N \pi_i \quad (3)$$

where N represents the number of TPs. A graph G , temporal partitioned in N subsets (TPs), is correct if:

- $\bigcup_{i=1}^N N(\pi_i) = \emptyset$: each node $v_i \in V$ is mapped to only one TP (here we do not consider cloning of operations in the DFG);
- $\bigcup_{i=1}^N N(\pi_i) = V$: all the nodes of V are mapped;
- $\forall \pi_i \in \wp, R(\pi_i) \leq R_{MAX}$: each TP fits in the resources available on the device;
- $\forall e_{ij} \in E, \pi(v_i) \geq \pi(v_j)$: the order of the execution of the TPs does not violate the dependencies among operations of the DFG (necessary condition to obtain the same functionality).

A correct set of TPs guarantees the same overall behavior of the original graph (when executed from 1 to N and considering a correct communication mechanism to transfer data among TPs). However, we are also interested on the minimization of the overall execution latency. The cost that reflects the overall execution latency in a time-multiplexed device can be estimated by the equation (1) or (2), when partial or full reconfiguration of the available resources is considered respectively. $CS(\wp)$ returns the minimum execution latency (number of control steps or clock cycles) of the partitioned solution, $CS(\pi_i)$ refers to the minimum execution latency of the TP π_i (it may include the communication costs and represents the execution latency of the critical path of the graph formed by the

subset of nodes in π_i and the correspondent edges, considering that nodes sharing FU instances can exist). ∂_i and ∂ represent the number of clock cycles to reconfigure the TP π_i or all the available resources respectively.

$$CS(\wp) = \sum_{i=1}^N CS(\pi_i) + \partial_i \quad (1)$$

$$CS(\wp) = \sum_{i=1}^N CS(\pi_i) + N \times \partial \quad (2)$$

The objective of our algorithm is to furnish a set of datapaths that will be executed in sequence with a minimum number of control steps². Each datapath unit fits on the physically available resources. For the sake of minimizing the number of TPs needed, exploiting sharing of FUs while doing temporal partitioning needs to be considered by the algorithm. Specifically, our algorithm has to output:

- The set of TPs (\wp): each TP identifying the nodes of the DFG assigned to it;
- The set of instances for each FU used (Φ);
- Each node of the DFG has to identify a specific FU instance of Φ implementing the operation.

From those outputs, it is straightforward to generate a behavioral HDL-RTL (hardware description language at the register transfer level) description of each TP control unit and a structural HDL-RTL description of each datapath, considering the existence of a HDL description for each FU. The configurations can be generated from those netlists using a traditional FPGA design flow.

3 Algorithm Simultaneously Exploiting Temporal Partitioning and Sharing of FUs

The algorithm uses an initial number of TPs that can be specified by the user. Another possibility is to use the number of levels of the DFG or the number of TPs utilized by any temporal partitioning algorithm without using sharing of FUs (e.g., ASAP [11]) as the initial number of TPs. The user has to specify the total number of available resources on the device. In addition, for each FU there exists a boolean variable which value indicates if the FU can be shared or not (sharing of some FUs may need more resources than the utilization of several FU instances, due to the overhead of using auxiliary circuits needed for the implementation of the sharing mechanism).

To a clear description, we show the main steps of the algorithm with a connection to Example 1. A brief exposi-

² We assume that each control/time step for scheduling is equal to the clock period of the system. Thus, there is no distinction among the use of clock cycle, control step or time step.

tion of the steps performed, when considering sharing of all FUs, is stretched in Figure 3.

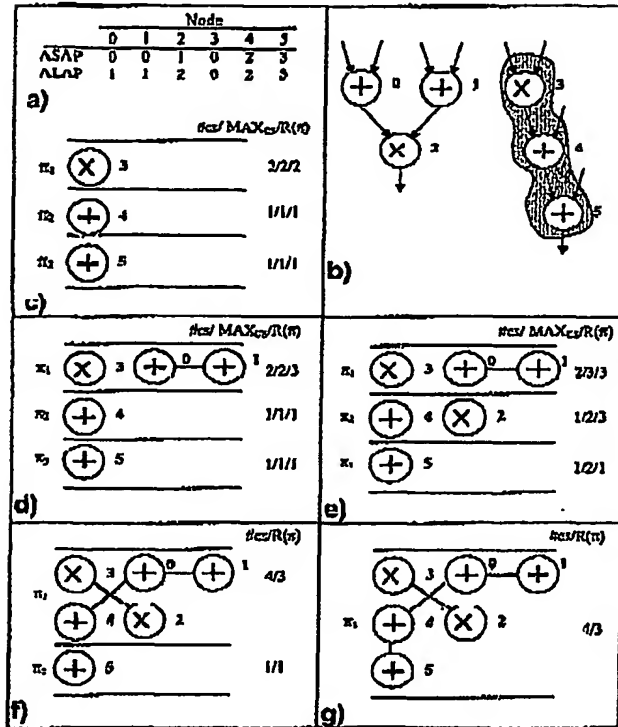


Figure 3. Algorithm execution through an example: a) ASAP and ALAP start times; b) The nodes in the critical path identified by the gray region; c), d), e), f) and g) show iterations of the algorithm.

The algorithm starts with the following steps:

1. Compute the set of nodes $child^3$ of each node of the DFG.
2. Map an FU instance to each operation in the DFG (at the moment neither consider more than one FU for the same operation nor FUs capable to implement more than one operation);
3. Estimate the area and execution latency of each node in the DFG according to the FU characterization, existent in the component library, for the target device. This step is beyond the scope of this article and from now on we will assume that there exists, for each FU, an estimation of the number of resources and of the execution latency;

³ A node v_i is child of a node v_j if there exists a path from v_j to the end of the DFG that includes v_i .

4. Perform the ASAP (as soon as possible) and ALAP (as late as possible) start times for each node in the DFG (see Figure 3a)), both unconstrained. When doing the ALAP scheme, the algorithm also calculates the ALAP level of each node;
5. Determine the set of nodes in one of the critical paths of the DFG (see Figure 3b));
6. Create a number of TPs equal to the input number specified (see the three TPs initially created in Figure 3c));
7. Assign each node of the set of nodes in one of the critical paths of the DFG (determined in point 5) to a TP by ascending level. When the number of TPs is larger than the number of nodes in the critical path, the last TPs are left empty; otherwise the last nodes of the set are left unassigned (see the nodes assigned to each TP in Figure 3c));
8. Assign the size (number of resources used) of a node in a TP to the current size of that TP (see Figure 3c)).

After the above steps the main kernel of the algorithm is executed (see the pseudo-code in Figure 4, Figure 5 and Figure 6). Some of the most important functions used by the algorithm are listed and briefly explained below:

- $v_i.ALAP_{level}()$: returns the level of v_i considering an ALAP leveling scheme;
- $v_i.ALAP_{start}()$: returns the ALAP start time of v_i ;
- $\pi_k.addEl(v_i)$: adds the node v_i to π_k ;
- $\pi_k.rmEl(v_i)$: removes v_i from π_k ;
- $\pi_k.sched(v_i)$: returns the number of control steps of the critical path considering that v_i is mapped to π_k ;
- $\wp.add(\pi_k)$: adds a new TP π_k to the current set of TPs (π_k will be the last TP in the set);
- $\wp.elAt(i)$: returns the i^{th} TP from the set of TPs (\wp);
- $findNodes(i)$: returns a list of nodes ready to be mapped to the i^{th} TP.

Our algorithm will be progressively constructing a global solution. On each iteration, the algorithm traverses the sequence of the existent TPs trying to assign ready nodes to each TP. Each TP has an associated maximum slot time (MAX_{CS}). A node ready to be mapped to a TP is only really considered for mapping if the resultant execution latency of that TP (considering the mapping) does not exceed the correspondent MAX_{CS} (line 15 of Figure 4 and lines 2, 21 and 29 of Figure 5). MAX_{CS} of a given TP π_i is equal to the critical path latency of that TP added by a relax amount: $CS(\pi_i) + relax$. On each iteration over the TPs the relax value is incremented by the great common divi-

sor (gcd) among all the execution latencies of the operations in the DFG (line 24 of Figure 4). When a node is mapped (see function *mapNode* in Figure 6), the critical path length of the associated TP is updated (lines 4 and 5 of Figure 6).

The algorithm considers that nodes in contiguous time steps mapped to the same TP and with the same operation should be bound to the same FU instance.

A list of nodes ready for mapping to a current TP is used. The list has the nodes sorted by increasing ALAP start times (the candidate operation having the least ALAP value will have the highest priority) and, for nodes with the same ALAP start time, it uses the ASAP start time as a tiebreak (by ascending or descending order). The list is determined examining for a given node its predecessors (they already must be mapped in TPs before the current TP) and the child set (the nodes child of the node to be mapped must be on TPs after the TP under consideration). The incremental update of the list of the nodes candidate to be mapped to the current TP, when each node is mapped, is an option of the algorithm (lines 6 and 7 in Figure 6). When such option is disabled the algorithm only tries to do update when the list is empty. The algorithm uses a static-based approach in the sense that the ALAP/ASAP values are calculated only once and they are no more time updated.

```

1. // begin main kernel
2. BitSet NodesSched = marked with the nodes already mapped to TPs;
3. int NumTP = 0; relax = 0;
4. int step = gcd(All nodes in DFG);
5. while(notAllNodesSched(NodesSched)) {
6.   Loop B: while(SchedNum < maxPartitions) {
7.     Vector listReady = findNodes(NumTP);
8.     Vector  $\pi_k$  = g.elAt(NumTP);
9.     while(!listReady.isEmpty()) {
10.      Node  $v_k$  = listReady.rmFirst();
11.      int  $R_{kpw}$  = R( $\pi_k$ ) + R( $v_k$ );
12.      Boolean fit = ( $R_{kpw}$  <=  $R_{max}$ );
13.      // CS( $\pi_k$ ) when  $v_k$  is mapped to  $\pi_k$ ;
14.      int  $CS_{new}$  =  $\pi_k$ .sched( $v_k$ );
15.      if(( $CS_{new}$  > ( $CS(\pi_k) + relax$ )) && (( $\pi_k$  is the last TP) && ( $N(\pi_k) \neq \emptyset$ ) || ( $v_k$ .ALAPlevel() <  $R(v_k)$ ))) {
16.        tryToSched( $R_{kpw}$ ,  $v_k$ , fit,  $CS_{new}$  -  $CS(\pi_k)$ ,  $\pi_k$ , NodesSched, update,  $CS_{new}$ , ListReady); Figure 5
17.      } else {
18.        tryToSched( $R_{kpw}$ ,  $v_k$ , fit, relax,  $\pi_k$ , NodesSched, update,  $CS_{new}$ , ListReady); Figure 5
19.      }
20.    }
21.    NumTP++;
22.  }
23.  NumTP = 0;
24.  relax += step;
25. }
26. // end main kernel

```

Figure 4. Main kernel of the proposed algorithm.

```

1. tryToSched(int  $R_{kpw}$ , Node  $v_k$ , Boolean fit, int relax, TP  $\pi_k$ , BitSet NodesSched, Boolean update, int  $CS_{new}$ , ListReady) {
2.   if(( $CS_{new}$  <= ( $relax + CS(\pi_k)$ )) || ( $CS(\pi_k) == 0$ )) {
3.     if only one node  $v_j$  in  $\pi_k$  {
4.       if( $v_j$ .ALAPstart() <  $v_k$ .ALAPstart()) {
5.         if(( $R_{kpw}$  - R( $v_j$ )) <=  $R_{max}$ ) {
6.            $\pi_k$ .rmEl( $v_j$ );
7.           R( $\pi_k$ ) = R( $v_k$ );
8.            $\pi_k$ .addEl( $v_k$ );
9.           NodesSched.clear( $v_j$ );
10.          NodesSched.set( $v_k$ );
11.          continue LOOP B;
12.        }
13.      }
14.      boolean canShare = try sharing with a node of the same type with a path of shared FUs with the smallest length (number of nodes);
15.      if(canShare) {
16.        if(fit && share produces increase) {
17.          canShare = false;
18.          rmShare( $v_k$ );
19.        } else {
20.          int  $CS_{new}$  =  $\pi_k$ .sched( $v_k$ );
21.          if( $CS_{new}$  < ( $relax + CS(\pi_k)$ )) {
22.            mapNode( $\pi_k$ ,  $v_k$ , NodesSched, update,  $CS_{new}$ , ListReady); Figure 6
23.          } else {
24.            rmShare( $v_k$ );
25.            canShare = false;
26.          }
27.        }
28.      }
29.      if(!canShare && fit && ( $CS_{new}$  <= ( $relax + CS(\pi_k)$ )) || ( $CS(\pi_k) == 0$ )) {
30.        mapNode( $\pi_k$ ,  $v_k$ , NodesSched, update,  $CS_{new}$ , ListReady); Figure 6
31.      }
32.      if( $v_k$  not mapped and no FU with operation type of  $v_k$  in thisTP and  $v_k$  does not fit and thisTP is the last TP) {
33.        create a new TP  $\pi_k$ ;
34.        g.add( $\pi_k$ );
35.        mapNode( $\pi_k$ ,  $v_k$ , NodesSched, update,  $CS(\pi_k)$ , ListReady); Figure 6
36.        break LOOP B;
37.      }
38.    }
39. } // end tryToSched

```

Figure 5. Function tryToSched.

```

1. mapNode(TP  $\pi_k$ , Node  $v_k$ , BitSet NodesSched, Boolean update, int  $CS_{new}$ , Vector ListReady) {
2.    $\pi_k$ .addEl( $v_k$ );
3.   NodesSched.set( $v_k$ );
4.   if( $CS_{new}$  >  $CS(\pi_k)$ )
5.      $CS(\pi_k)$  =  $CS_{new}$ ;
6.   if(update)
7.     updateAndSortALAP(ListReady,  $v_k$ );
8. } // end mapNode

```

Figure 6. Function mapNode.

A special directed edge between two nodes is used to identify in the DFG that both nodes share the same FU instance. A path of nodes connected by such edges identifies the sequence of utilization of the FU instance (from the source to the sink). Those edges are added by the algorithm to the initial DFG and provide an efficient way to both represent sharing of FU instances, to determine the execution latency and to generate the datapath and control unit descriptions. When a node v_n is bound to an FU instance that is already shared by two or more operations, the algorithm adds a special edge from the last node of the path of shared FUs to v_n . However, when there are one or more nodes in that path that are child of v_n in the DFG, v_n is inserted immediately before the first child found traversing the path from left to right (see Figure 7).

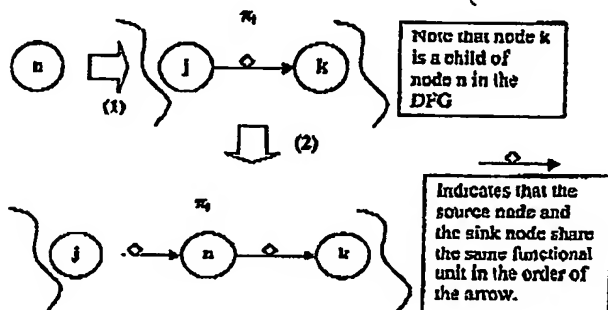


Figure 7. Special mapping of a node in a TP.

The algorithm considers the interchange⁴ of a node previously assigned to a TP with a node ready to be mapped in that TP. This occurs if the TP has only one node, the node to be mapped has lower ALAP start time and the change is feasible in terms of available resources and MAX_{CS} (lines from 3 to 13 of Figure 5).

After the execution of the main kernel the algorithm considers a merge operation. Such operation tries to group adjacent TPs in a single TP considering resource sharing among FU instances of the same type in two considered TPs (see Figure 3e), f) and g)). This step of the algorithm is done incrementally and it stops when no merging is possible. From the two first TPs of Figure 3e) it can be seen that sharing the FUs, the two TPs can be merged. Figure 3f) illustrates the result of that merge. From the later result we can figured out that π_1 and π_2 are also able to be merged (with node 5 sharing the FU instance already shared by nodes 0, 1 and 4). The resultant merge is shown in Figure 3g) and the final solution requires only a single TP.

When a node does not fit (in the sense that the MAX_{CS} of the considered TP is violated) in the last TP and if that

⁴ This involves the move of a node from a TP to the list and the map of the node ready to that TP.

TP is empty then the algorithm prefers to relax the TP to accommodate the node (see lines 15 to 17 in Figure 4).

When finding if a node to be mapped can share an FU instance already existent in the considered TP, the algorithm binds the node/operation to the FU instance with lower length path of nodes sharing it (line 14 in Figure 5).

The steps of the algorithm shown in lines 32 to 37 of Figure 5 are related to the case that when a node does not fit (with or without resource sharing or with relaxation of MAX_{CS}) in the last TP a new TP is added to the existent set of TPs and the node is mapped to it.

4 Experimental results

All the algorithms considered in this section have been implemented with the JavaTM language. All the executions of the algorithms were conducted on a portable PC (Pentium-II @366MHz, 196Mb RAM) with the JIT compiler integrated in the JDK1.2 running on Windows98.

The algorithm proposed in this paper has been tested with a number of representative examples (with variable complexity) and, whenever possible, the results obtained are compared with other approaches.

The SEWHA is the auto regression filter presented in [13], HAL is the loop body of the differential equation example [14], EWF is a digital fifth order elliptic wave filter [15], FIR is a 12-tap finite impulse filter, and Mat4x4 corresponds to a fully parallel multiplication of two matrices with 4x4 integer elements each one. Mat4x4 is used as an example of high operation level parallelism degree.

We consider for all the experiments an execution latency of 1 clock cycle for each adder and 2 clock cycles for each multiplier. For the number of resources needed for each FU we consider 1 cell for each adder and 4 cells for each multiplier. Table II shows the main characteristics of the considered examples (number of operations of each type, total number of resources and the critical path length).

Table II. Characteristics of the examples.

Example	Number of operations (\times , $+/$)	Total number of operations	Total number of resources (cells)	Critical path length of the DFG (clock cycles)
Ex1	(2, 4)	6	12	4
EWF	(8, 26)	34	58	17
FIR	(12, 11)	33	59	6
HAL	(4, 6)	10	22	7
SEWHA	(16, 12)	28	76	11
Mat4x4	(64, 48)	112	304	4

4.1 Sharing versus not sharing

Table III shows results for the considered examples. Our* and Our** identify results obtained by applying the proposed algorithm. Our* considers resource sharing for both adder and multiplier units, and Our** only considers resource sharing for multiplier units. #cs identifies the execution latency (number of clock cycles) and #p the number of TPs. Each solution related to our algorithm was obtained in less than 1s of CPU time.

Table III. Results obtained for the examples.

Example	R _{MAX}	Approach							
		ASAP		SA		Our*		Our**	
		#p	#cs	#p	#cs	#p	#cs	#p	#cs
ExI	6	2	5	2	4	1	4	2	4
SEHWA	6	18	36	18	35	1	34	6	34
	10	9	24	9	19	1	18	6	18
	15	6	18	6	15	1	15	5	15
HAL	6	5	11	5	9	2	10	3	10
	10	3	10	3	7	2	7	3	7
EWF	6	12	26	12	23	1	23	9	25
	10	6	22	6	22	5	18	7	18
	15	4	19	4	18	1	17	4	18
FIR	6	14	28	14	27	1	20	4	27
	10	7	16	7	15	1	15	4	15
	15	5	12	5	11	1	11	3	11
MAT4x4	6	72	136	72	134	1	130	21	130
	10	37	69	37	69	1	66	17	66
	15	25	47	25	46	1	46	10	46
	20	16	29	16	29	2	29	4	29

The SA results were obtained with a simulated annealing version to do temporal partitioning without resource sharing proposed in [16]. Here, the algorithm is tuned to optimizing the overall execution time (the algorithm can also exploit the tradeoff between execution time and communication costs). The ASAP results refer to the leveling technique proposed in [11].

Only Mat4x4 needed to start with the number of TPs obtained by the ASAP approach to achieve the best solution. For all the other examples, the best solution was obtained starting with an initial number of TPs equal to the number of levels of the DFG. The results for Mat4x4 in Table III were collected disabling the update of the list of nodes ready for each node mapped (the list is updated only when it is empty). It is strongly recommended to disable the update option for examples with high-level degree of parallelism and a small critical path length.

The values in bold in the 6th and 8th columns of Table III show the minimum execution latency for the datapaths obtained by the considered approaches (not considering configuration times). The values in bold in the 10th column represent that, even without considering sharing of adders, our algorithm returns solutions with execution latencies

equal to the execution latencies obtained sharing all the resources (8th column), despite the fact that those solutions need more TPs.

When considering resource sharing for all FUs, a minimum number of TPs (only 4 cases of Table III needed more than one TP to produce a minimum execution time) seems to ensure solutions with lower execution latencies than the obtained by doing temporal partitioning with ASAP or SA for the majority of the examples (only one case is not as good as SA). Note that when all the FUs can be shared and the resource overhead to implement sharing is not taken into account, an empirical observation tell us that the solutions with lower execution latency are those with only one TP. This is expected by the fact that a new TP produces an equal or worse effect than sharing FU instances on the overall execution latencies because all the nodes in that TP can only start executing after the end of the execution of the TP immediately before.

When sharing of adders is not considered the algorithm is capable to find 13 solutions without inferior execution latency.

4.2 Exploiting the number of TPs

An exploitation of the overall execution latency versus the number of TPs is shown in Figure 8. Those results were produced by calling the algorithm several times, each time starting with a different initial number of TPs from a range of 1 to 15. The exploitation has been done in approximately 5.4s of CPU time. All the solutions use only a single TP and the best result (execution latency equal to 66 clock cycles) has been achieved when the algorithm started with 8 TPs. The results without considering sharing of adders are shown in Figure 9. The algorithm exploited a range of TPs from 1 to 26 and the minimum execution latency achieved was 66 clock cycles (solution with 21 TPs). Based on those results we can select a solution that minimizes the global execution latency taking into account the reconfiguration times (see equation (2)).

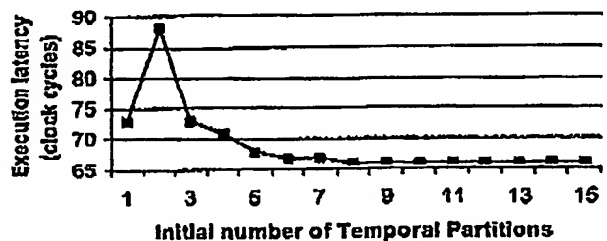


Figure 8. Execution latency versus the initial number of TPs for Mult4x4 obtained by the proposed algorithm, when R_{MAX}=10 (sharing of adders and multipliers).

From the results presented so far we may conclude that sharing FUs can reduce the number of TPs without increasing the overall execution time. Moreover, a minimum number of TPs can be a priority, when an FPGA with significant reconfiguration times is used. Due to its low computational complexity, the algorithm can be used to exploit the design space based on the tradeoff between the number of TPs and the overall execution latency.

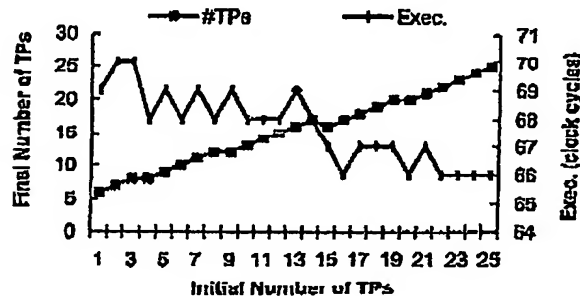


Figure 9. Execution latency and the final number of TPs obtained by the algorithm for Mult4x4, when $R_{MAX}=10$ (no sharing of adders).

4.3 Comparison with other schedulers

At this point a question may occur: is the algorithm competitive when a single TP is envisaged? Table IV shows results for EWF and SEHWA, considering various sizes for the available resources (R_{MAX}). The schedules obtained by the proposed algorithm considering only one TP are shown (see the 5th column). The number of resources used for each type of FU for each solution is also shown (last column). "Fixed" refers to results collected from the state-of-the-art schedulers [17][18][19] and represent optimal (identified with *) or near-optimal scheduling results (without enter into account with temporal partitioning) for the specified constraint on the number of FUs for each type of operation (see the 2nd column). The results show that our algorithm is efficient, even when we are interested on a final solution with a single TP.

The result labeled with a "*" is achieved without an incremental update of the list of the nodes ready to be mapped. This result shows that the algorithm did not skip from a local minimum, since at least the result related to $R_{MAX}=15$ should be achieved. The first 4 results obtained for SEHWA consider the increasing order of the ASAP values as the second key (there is no evidence to suggest when it is better to use the decreasing or the increasing ASAP values as the second key).

The number of each FU instance allocated by our algorithm for each R_{MAX} constraint only was different in two

cases to the constraints used (with total number of resources equal to R_{MAX}) to produce the near-optimal scheduling results (see Table IV). Therefore, it seems that our algorithm can also be used to a fast identification of the number of FU instances needed, considering a specific number of maximum resources available on the device.

Table IV. Comparison of scheduling results obtained for EWF and SEHWA.

Example	Approach				
	Fixed [17][18][19]		Our		
	constraints (x, +)	#cs	constraints R_{MAX} (#p=1)	#cs	(x, +)
EWF	(1, 2)	21*	6	23	(1, 2)
	(1, 3)	18	7	22	(1, 3)
	(2, 1)	21*	9	22	(1, 5)
	(2, 2)	18*	10	20	(2, 2)
	(2, 3)	18	11	18	(2, 3)
	(3, 2)	18	14	18	(2, 6)
SEHWA	(3, 3)	17*	15	17	(3, 3)
	(1, 1)	34*	5	34	(1, 1)
	(2, 1)	18*	9	20	(2, 1)
	(2, 2)	18*	10	18	(2, 2)
	(3, 1)	16*	13	17	(3, 1)
	(3, 2)	15*	14	15	(3, 2)
	(3, 3)	15*	15	15	(3, 3)
	(4, 1)	16*	17	16*	(4, 1)
	(4, 2)	11*	18	11	(4, 2)

5 Related Work

As far as we know, the development of temporal partitioning algorithms was firstly considered in [9][2]. The similarities of both scheduling on high-level synthesis [8] and temporal partitioning allow the use of common scheduling schemes for partitioning. Some authors, such as [9][10], have considered temporal partitioning at behavioral levels having in mind the integration of synthesis.

In [9], a heuristic based on a static list scheduling algorithm, enhanced to consider temporal partitioning and partial reconfiguration, is shown. The approach exploits the dynamic reconfiguration capability of the devices, while doing temporal partitioning.

In [10][20] the temporal partitioning problem is modeled in a specified 0-1 non-linear programming (NLP) model. The problem is transformed to integer linear programming (ILP) and the solution determined by an ILP solver. Due to the long execution times, this approach is not practical for large input examples. Some heuristic methods have been developed to permit its usability on larger input examples [21]. Kaul [22] exploits the loop fission technique while doing temporal partitioning in the presence of loops to minimize the overall latency by utilization of the active TP as long as possible. Sharing of

functional units is considered inside tasks and temporal partitioning is performed at the task level. Design space exploitation is performed by inputting to the temporal partitioning algorithm different design solutions for each task. Such solutions are generated by a high-level synthesis tool (constraining the number of FUs of each type). This approach lacks a global view and is time-consuming.

The simplest approaches only consider temporal partitioning without exploiting sharing of FUs. In [11], both a temporal partitioning algorithm based on leveling the operations by an ASAP scheme and other based on clustering a number of nodes are used. The algorithm fills the available resources in the increasing order of the ASAP levels. The selection of nodes in the same level is arbitrary and the algorithm switches to another TP when it encounters the first node that does not fit on the current TP. The approach does not consider neither communications costs nor resource sharing. In [23] another algorithm is presented that selects the nodes to be mapped in a TP with two different approaches (one for satisfying parallelism and an other for decreasing communication costs). In [12], an algorithm based on the extension of the ASAP or ALAP leveling schemes resorting to the mobility of each node to select among the nodes has been considered. [12] also shows an algorithm that searches recursively in the list of ready nodes so that if a node cannot be mapped to the current partition, other nodes can be considered.

[16] considers both communication costs among different TPs that can occur and the overall execution time. The authors presented an extension to static list scheduling, which permits to the algorithm sensitivity to the communication costs while trying to minimize the overall execution time. The results presented, when compared to near-optimal solutions obtained with a simulated annealing algorithm tuned to do temporal partitioning while minimizing an objective function, that integrates the execution time of the TPs and the communication costs, revealed the efficiency of the approach.

[24] presents a method to do temporal partitioning considering pipelining of the reconfiguration and execution stages. The approach divides an FPGA into two portions to overlap the execution of a TP in one portion (previously reconfigured) with the reconfiguration of the other portion.

In [25] constraint logic programming is used to solve temporal partitioning, scheduling, and dynamic module allocation. However, the approach needs a specification of the number of each FU before processing and may suffer of long runtimes.

More related to our approach is the algorithm presented in [26]. A scheme based on the force-directed list scheduling algorithm that considers resource sharing and temporal partitioning is shown. The algorithm tries to minimize the overall execution time, performing a tradeoff between the number of TPs and sharing of FUs. However, the approach adapted a scheduling algorithm not originally tailored to

do temporal partitioning and lacks of a global view. Instead, our approach proposes a novel algorithm matched to the combination of temporal partitioning and sharing of FUs that maintains a global view.

6 Conclusions and Future Work

In this paper we have presented a new and useful algorithm combining temporal partitioning, sharing of functional units, scheduling, allocation and binding. Unlike other approaches, this algorithm merges those tasks in a combined and global method. The obtained results, from a number of benchmarks, strongly confirm the efficiency and effectiveness of the idea.

The low computation time achieved, when dealing with the presented examples, shows that the algorithm is fast and efficient and thus can be used on large examples.

The inclusion of functional units with pipeline stages and the consideration of more than one implementation for a given operation will be considered in a near future. Another important issue is the overlapping of reconfiguration and execution that should be considered by future enhancements. Finally, aspects related to conditional paths and loops will also need to be focused of future work.

7 Acknowledgments

The author would like to acknowledge the support from the "Fundação para a Ciência e Tecnologia" and particularly from the PRAXIS XXI Program under the scope of the AXEL Project (PRAXIS/P/EEI/12154/1998).

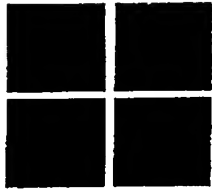
8 References

- [1] A. DeHon, *Reconfigurable Architectures for General Purpose Computing*, PhD Thesis, AI Technical Report 1586, MIT, 545 Technology Sq., Cambridge, MA02139, Sept. 1996, <http://www.ai.mit.edu/people/andre/phd.html>.
- [2] X.-P. Long, H. Amano, "WASMII: a Data Driven Computer on a Virtual Hardware," in *Proc. of the 1st IEEE Workshop on Field Programmable Custom Computing Machines*, Napa Valley, CA, USA, April 5-7, 1993, pp. 33-42.
- [3] Xilinx Inc., *Virtex Field Programmable Gate Arrays*, version 1999.
- [4] R. D. Hudson, D. I. Lehn, and P. M. Athanas, "A Run-Time Reconfigurable Engine for Image Interpolation," in *Proc. of the 6th IEEE Symposium on Field-Programmable Custom Computing Machines*, Napa Valley, CA, USA, April 15-17, 1998.
- [5] T. Fujii, et al, "A Dynamically Reconfigurable Logic Engine with a Multi-Context/Multi-Mode Unified-Cell Architecture," in *Proc. of the IEEE Int'l Solid State Circuits Conference*, SA, CA, Feb. 15-17, 1999. See

<http://www.nec.co.jp/english/today/newsrel/9902/1502.html>

- [6] S. Trimberger, "Scheduling Designs into a Time-Multiplexed FPGA," in *Proc. of the 6th ACM Int'l Symposium on FPGAs*, Monterey, CA, USA, February 22-24, 1998.
- [7] H. Liu and D. F. Wong, "Circuit partitioning for dynamically reconfigurable FPGAs," in *Proc. 7th ACM Int'l Symposium on FPGAs*, Monterey, CA, USA, Feb. 21-23, 1999, pp. 187-194.
- [8] D. Gajski, et al., *High-Level Synthesis, Introduction to Chip and System Design*, Kluwer Academic Publishers, 1992.
- [9] M. Vasilko, D. Ait-Boudaoud, "Architectural Synthesis Techniques for Dynamically Reconfigurable Logic," in *Proc. of the 6th Int. Workshop on Field-Programmable Logic and Applications*, Darmstadt, Germany, Sept. 23-25, 1996, LNCS, vol. 1142, Springer-Verlag, pp. 290-296.
- [10] I. Ouaiss, et al., "An Integrated Partitioning and Synthesis System for Dynamically Reconfigurable Multi-FPGA Architectures," in *Proc. of the 5th Reconfigurable Architectures Workshop*, Orlando, Florida, USA, March 30, 1998, pp. 31-36.
- [11] Karthikeya M. GajjalaPurna, and Dinesh Bhatia, "Temporal Partitioning and Scheduling Data Flow Graphs for Reconfigurable Computers," in *IEEE Transactions on Computers*, vol. 48, no. 6, June 1999, pp. 579-591.
- [12] João M. P. Cardoso, and Horácio C. Neto, "Macro-Based Hardware Compilation of Java™ Bytecodes into a Dynamic Reconfigurable Computing System," in *Proc. of the IEEE 7th Symposium on Field-Programmable Custom Computing Machines (FCCM'99)*, Napa Valley, CA, USA, April 21-23, 1999, Kenneth L. Pocek and Jeffrey Arnolds (eds.), IEEE Computer Society Press, Los Alamitos, CA, USA, pp. 2-11.
- [13] R. Jain, A. Parker, N. Park, "Module Selection for Pipelined Synthesis," in *Proc. of the 25th Design Automation Conference*, Anaheim, CA, USA, June 12-15, 1988, pp. 542-547.
- [14] P. G. Paulin, J. P. Knight, and E. F. Girczyc, "HAL: A Multi-Paradigm Approach to Automatic Data Path Synthesis," in *Proc. of the 23rd Design Automation Conference*, Las Vegas, NV, USA, June 29-July 2, 1986, 263-270.
- [15] P. Dewilde, E. Deprettere, and R. Nouta, "Parallel and Pipelined VLSI Implementation of Signal Processing Algorithms," in *VLSI and Modern Signal Processing*, S.Y. Kung, H.J. Whitehouse, T.Kailath (eds.), Prentice-Hall 1985, pp. 258-264.
- [16] João M. P. Cardoso, and Horácio C. Neto, "An Enhanced Static-List Scheduling Algorithm for Temporal Partitioning onto RPLs," in *Proc. of the IFIP X International Conference on Very Large Scale Integration (VLSI'99)*, Lisbon, December 1-3, 1999, Luis M. Silveira, Srinivas Devadas and Ricardo Reis (eds.), *VLSI: Systems on a Chip*, Kluwer Academic Publishers, pp. 485-496.
- [17] M. Narasimhan, and J. Ramanujam, "A Fast Approach to Computing Exact Solutions to the Resource-Constrained Scheduling Problem," to appear in *ACM Transactions on Design Automation of Electronic Systems (TODAES)*, Vol. 7, No. 1, January 2002. URL: <http://www.acm.org/todaes/VTN1/TOC.html>
- [18] P.-Y. Hsiao, G.-M. Wu, and J.-Y. Su, "MPT-based branch-and-bound strategy for scheduling problem in high level synthesis," in *IEE Proc. Computers and Digital Techniques*, Vol. 145, No. 6, November 1998, pp. 425-432.
- [19] M. K. Dhodhi, F. H. Hiescher, K. M. Storer, and F. Bhasker, "Datapath Synthesis Using a Problem-Space Genetic Algorithm," in *IEEE Transactions on CAD of Integrated Circuits and Systems*, Vol. 14, No. 8, August 1995, pp. 934-944.
- [20] M. Kaul, R. Vemuri, "Optimal Temporal Partitioning and Synthesis for Reconfigurable Architectures," in *Proc. of the Design, Automation & Test in Europe*, Paris, France, Feb. 23-26, 1998, pp. 389-396.
- [21] M. Kaul, R. Vemuri, "Temporal Partitioning combined with Design Space Exploration for Latency Minimization of Run-Time Reconfigured Designs," in *Proc. of Design, Automation & Test in Europe*, Paris, France, Feb. 23-26, 1999, pp. 202-209.
- [22] Meenakshi Kaul, Ranga Vemuri, Sritam Govindarajan, Iyad E. Oualis, "An Automated Temporal Partitioning and Loop Fission approach for FPGA based reconfigurable synthesis of DSP applications," in *Proc. of the IEEE/ACM Design Automation Conference (DAC'99)*, New Orleans, LA, USA, June 21-25, 1999, pp. 616-622.
- [23] Arusshi Takayama, Yuichiro Shibata, Katsuko Iwai, and Hideharu Amano, "Dataflow Partitioning and Scheduling Algorithms for WASMII, a Virtual Hardware," in *Proc. of the 10th Int. Conference on Field-Programmable Logic and Applications (FPL'00)*, Villach, Austria, August 27-30, 2000, Reiner W. Hartenstein and Herbert Grünbacher (eds.), LNCS, vol. 1896, Springer-Verlag, Berlin, pp. 685-694.
- [24] Satish Ganesan, and Ranga Vemuri, "An Integrated Temporal Partitioning and Partial Reconfiguration Technique for Design Latency Improvement," in *Proc. of Design, Automation & Test in Europe (DATE'00)*, Paris, France, March 27-30, 2000, pp. 320-325.
- [25] Xue-jie Zhang, Kam-wing Ng, and Wayne Luk, "A Combined Approach to High-Level Synthesis for Dynamically Reconfigurable Systems," in *Proc. of the 10th Int. Conference on Field-Programmable Logic and Applications (FPL'00)*, Villach, Austria, August 27-30, 2000, Reiner W. Hartenstein and Herbert Grünbacher (eds.), LNCS, vol. 1896, Springer-Verlag, Berlin, pp. 361-370.
- [26] Awartika Pandey and Ranga Vemuri, "Combined Temporal Partitioning and Scheduling for Reconfigurable Architectures," in *Proc. SPIE Photonics East Conference, SPIE 3844*, Boston, Massachusetts, USA, Sept. 20-21, 1999, John Schewel, et al. (eds.), *Reconfigurable Technology: FPGAs for Computing and Applications*, pp. 93-103.

Confidential



PACT

XPP Technologies

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor

Patent

© 2002 PACT XPP Technologies

Version 0.1

2002-1-18

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 2

1 Introduction

This document describes a method for compiling a subset of a high-level programming language (HLL) like C or FORTRAN, extended by port access functions, to a reconfigurable data-flow processor (RDFP) as described in Section 3. I.e., the program is transformed to one or several configurations of the RDFP.

This method can be used as part of an extended compiler for a hybrid architecture consisting of standard host processor and a reconfigurable data-flow coprocessor. The extended compiler handles a full HLL like standard ANSI C. It maps suitable program parts like inner loops to the coprocessor and the rest of the program to the host processor. However, this extended compiler is not subject of this document.

2 Compilation Flow

This section briefly describes the phases of the compilation method.

2.1 Frontend

The compiler uses a standard frontend which translates the input program (e.g. a C program) into an internal format (IF) consisting of an abstract syntax tree (AST) and symbol tables. The frontend also performs well-known compiler optimizations as constant propagation, dead code elimination, common subexpression elimination etc. For details, refer to any compiler construction textbook like [1]. E.g., the SUIF compiler [2] can be used for this purpose.

2.2 Temporal Partitioning

Next, the program's IF representation is partitioned into sections which are executed sequentially on the RDFP by separate configurations. If the entire program can be executed by one configuration (fitting on the given RDFP), no temporal partitioning is necessary. This phase generates reconfiguration statements which load and remove the configurations sequentially according to the original program's control flow.

2.3 Configuration Generation

Finally, the program sections determined by the temporal partitioning are mapped to RDFP configurations. This phase generates a program code or data structure which is then used to directly program the RDFP.

3 Configurable Objects and Functionality of a RDFP

This section describes the configurable objects and functionality of a RDFP. A possible implementation of the RDFP architecture is a PACT XPP™ Core. Here we only describe the minimum requirements for a RDFP for this compilation method to work. The only data types considered are multi-bit words called *data* and single-bit control signals called *events*. Data and events are always processed as *packets*, cf. Section 3.2.

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 3

3.1 Configurable Objects and Functions

An RDFP consists of an array of configurable objects and a communication network. Each object can be configured to perform certain functions (listed below). It performs the same function repeatedly until the configuration is changed. The array needs not be completely uniform, i. e. not all objects need to be able to perform all functions. E. g., a RAM function can be implemented by a specialized RAM object which cannot perform any other functions. It is also possible to combine several objects to a "macro" to realize certain functions. Several RAM objects can, e. g., be combined to realize a RAM function with larger storage.

After a configuration has been removed, all information is lost. Only the contents (values) of a RAM are preserved during reconfiguration.

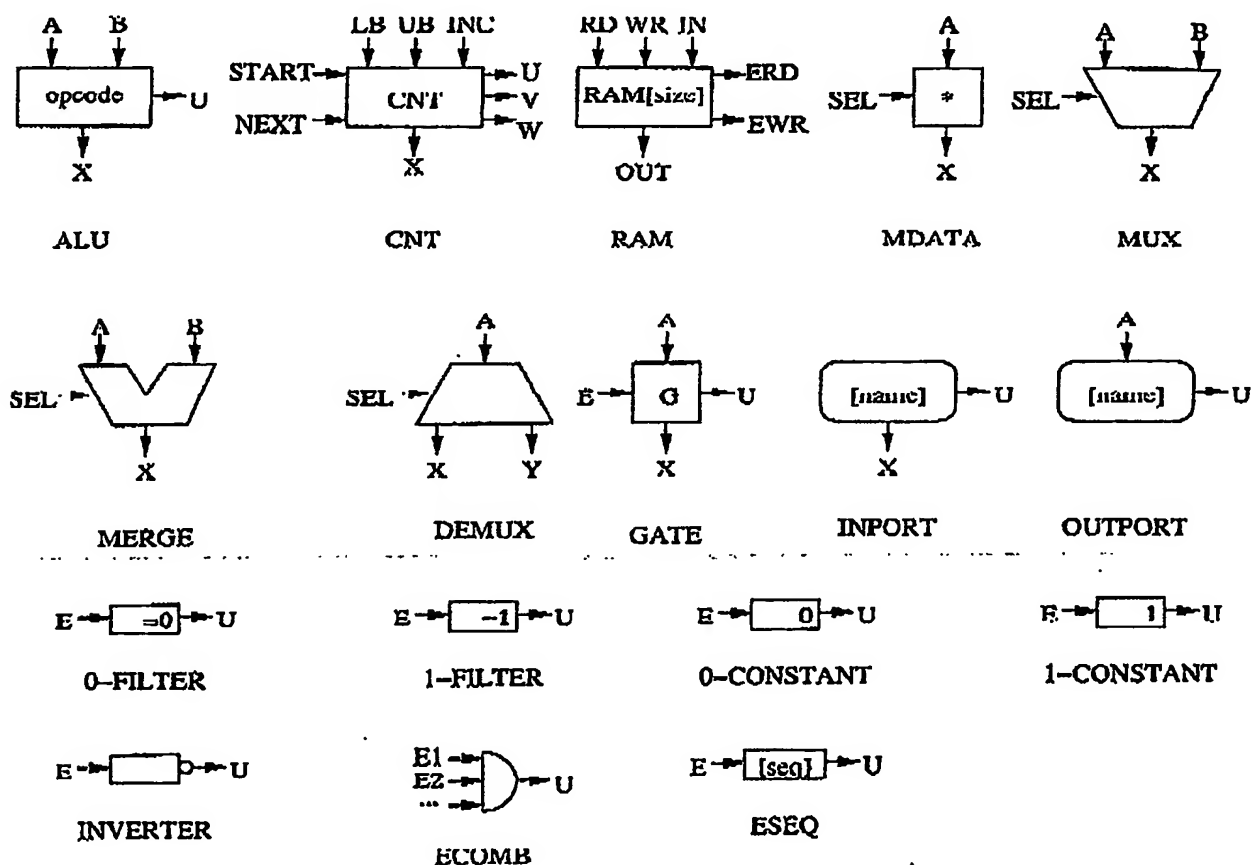


Figure 1: Functions of an RDFP

The following functions mainly handling data packets can be configured in an RDFP. See Fig. 1 for a graphical representation.

- **ALU[opcode]:** ALUs perform common arithmetical and logical operations on data. ALU functions ("opcodes") must be available for all operations used in the HLL.¹ ALU functions have two

¹Otherwise programs containing operations which do not have ALU opcodes in the RDFP must be excluded from the

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 4

data inputs A and B, and one data output X. Comparators have an event output U instead of the data output. They produce a 1-event if the comparison is true, and a 0-event otherwise.

- **CNT:** A counter function which has data inputs LB, UB and INC (lower bound, upper bound and increment) and data output X (counter value). A packet at event input START starts the counter, and event input NEXT causes the generation of the next output value (and output events) or causes the counter to terminate if UB is reached. If NEXT is not connected, the counter counts continuously. The output events U, V, and W have the following functionality: For a counter counting N times, N-1 event packets with value 0 (0-events) and one event packet with value 1 (1-event) are generated at output U. At output V, N 0-events are generated, and at output W, N 0-events and one 1-event are created. The 1-event at W is only created after the counter has terminated, i. e. a NEXT event packet was received after the last data packet was output.
- **RAM[size]:** The RAM function stores a fixed number of data words ("size"). It has a data input RD and a data output OUT for reading at address RD. Event output ERD signals completion of the read access. For a write access, data inputs WR and IN (address and value) and data output OUT is used. Event output EWR signals completion of the write access. ERD and EWR always generate 0-events. Note that external RAM can be handled as RAM functions exactly like internal RAM.
- **GATE:** A GATE synchronizes a data packet at input A back and an event packet at input E. When both have arrived, they are both inputs consumed. The data packet is copied to output X, and the event packet to output U.
- **MUX:** A MUX function has 2 data inputs A and B, an event input SEL, and a data output X. If SEL receives a 0-packet, input A is copied to output X and input B discarded. For a 1-packet, B is copied and A discarded.
- **MERGE:** A MERGE function has 2 data inputs A and B, an event input SEL, and a data output X. If SEL receives a 0-packet, input A is copied to output X, but input B is *not* discarded. The packet is left at the input B instead. For a 1-packet, B is copied and A left at the input.
- **DEMUX:** A DEMUX function has one data input A, an event input SEL, and two data outputs X and Y. If SEL receives a 0-packet, input A is copied to output X, and no packet is created at output Y. For a 1-packet, A is copied to Y, and no packet is created at output X.
- **MDATA:** A MDATA function multiplies data packets. It has a data input A, an event input SEL, and a data output X. If SEL receives a 1 packet, a data packet at A is consumed and copied to output X. For all subsequent 0-packets at SEL, a copy of the input data packet is produced at the output without consuming new packets at A. Only if another 1-packet arrives at SEL, the next data packet at A is consumed and copied.²
- **IMPORT[name]:** Receives data packets from outside the RDFP through input port "name" and copies them to data output X. If a packet was received, a 0-event is produced at event output U, too. (Note that this function can only be configured at special objects connected to external busses.)
- **EXPORT[name]:** Sends data packets received at data input A to the outside of the RDFP through output port "name". If a packet was sent, a 0-event is produced at event output U, too. (Note that this function can only be configured at special objects connected to external busses.)

supported HLL subset or substituted by "macros" of calling functions.

²Note that this can be implemented by a MERGE with special properties on XPP.

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 5

Additionally, the following functions manipulate only event packets:

- **0-FILTER, 1-FILTER:** A FILTER has an input E and an output U. A 0-FILTER copies a 0-event from E to U, but 1-EVENTs at E are discarded. A 1-FILTER copies 1-events and discards 0-events.
- **INVERTER:** Copies all events from input E to output U but inverts its value.
- **0-CONSTANT, 1-CONSTANT:** 0-CONSTANT copies all events from input E to output U, but changes them all to value 0. 1-CONSTANT changes all to value 1.
- **ECOMB:** Combines two or more inputs E1, E2, E3..., producing a packet at output U. The output is a 1-event iff one or more of the input packets are 1-events (logical *or*). A packet must be available at all inputs before an output packet is produced.³
- **ESEQ[seq]:** An ESEQ generates a sequence "seq" of events, e.g. "0001", at its output U. If it has an input START, one entire sequence is generated for each event packet arriving at U. The sequence is only repeated if the next event arrives at U. However, if START is not connected, ESEQ constantly repeats the sequence.

3.2 Packet-based Communication Network

The communication network of an RDFP can connect an outputs of one object (i. e. its respective function) to the input(s) of one or several other objects. This is usually achieved by busses and switches. By placing the functions properly on the objects, many functions can be connected arbitrarily up to a limit imposed by the device size. As mentioned above, all values are communicated as packets. A separate communication network exists for data and event packets. The packets synchronize the functions in a data-flow fashion. I. e., the function only executes when all input packets are available (apart from the exceptions where not all inputs are required as described above). The function also stalls if the last output packet has not been consumed. Therefore a data flow graph mapped to an RDFP self-synchronizes its execution without the need for external control. Only if two or more function-outputs are connected to the same function input (N to 1 connection), the self-synchronization is disabled. The user has to ensure that only one packet arrives at a time. Otherwise a packet might get lost, and the value resulting from combining two or more packets is undefined. Therefore this should be avoided. However, a function output can be connected to many function inputs (1 to N connection) without problems.

There are some special cases:

- A function input can be *preloaded* with a distinct value during configuration. This packet is consumed like a normal packets coming from another object.
- A function input can be defined as *constant*. In this case, the packet at the input is reproduced repeatedly for each function execution. It is even possible to connect an output of another function to a constant input. In this case, the constant value is changed as soon as a new packet arrives at the input. Note that there is no self-synchronization in this case, too. The function is not stalled until the new packet arrives since the old packet is still used and reproduced.

³Note that this function is implemented by the EAND operator on the XPP.

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 6

An RDFFP requires register delays in the dataflow. Otherwise very long combinational delays and asynchronous feedback is possible. We assume that delays are inserted at the inputs of some functions (like for most ALUs) and in some routing segments of the communication network.

4 Temporal Partitioning

The details of Temporal Partitioning need to be inserted from J. Cardoso's documents.

5 Configuration Generation

5.1 Language Definition

The following HLL features are not supported by the method described here:

- pointer operations
- library calls, operating system calls (including standard I/O functions)
- recursive function calls (Note that non-recursive function calls can be eliminated by function inlining and therefore are not considered here.)
- All scalar data types are converted to type integer. Integer values are equivalent to *data packets* in the RDFFP. Arrays (possibly multi-dimensional) are the only composite data types considered.

The following additional features are supported:

INPORTS and OUTPORTS can be accessed by the HLL functions *getstream(name, value)* and *putstream(name, value)* respectively.

5.2 Mapping of High-Level Language Constructs

This method converts a HLL program to a control/data-flow graph (CDFG) consisting of the RDFFP functions defined in Section 3.1. Before the processing starts, all HLL program arrays are mapped to RDFFP RAM functions. An array *x* is mapped to RAM *RAM(x)*. If several arrays are mapped to the same RAM, an offset is assigned, too. The RAMs are added to an initially empty CDFG. There must be enough RAMs of sufficient size for all program arrays.

The CDFG is generated by a traversal of the AST of the HLL program. The following two pieces of information are maintained at every program point⁴ during the traversal:

⁴In a program, *program points* are between two statements or before the beginning or after the end of a program structure like a loop or a conditional statement.

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 7

- START points to an event output of an object. It delivers a 0-event whenever the program execution at this program point starts. At the beginning, a 0-CONSTANT preloaded with an event input is added to the CDFG. (It delivers a 0-event immediately after configuration.) START initially points to its output. The STOP signal generated after a program part has finished executing is used as new START signal for the next program part or signals termination of the entire program.
- VARLIST is a list of {variable, object-output} pairs. The pairs map integer variables (no arrays) to a CDFG object's output. The first pair for a variable in VARLIST contains the output of the object which produces the value of this variable valid at the program point. New pairs are always added to the front of VARLIST. The expression VARDEF(var) refers to the the object-output of the first pair with variable var in VARLIST.⁵

The following subsections systematically list all ILL components and describe how they are processed, thereby altering the CDFG, START and VARLIST.

5.2.1 Integer Expressions and Assignments

Straight-line code without array accesses can be directly mapped to a data-flow graph. One ALU is allocated for each operator in the program. Because of the self-synchronization of the ALUs, no explicit control or scheduling is needed. Therefore processing these assignments does not access or alter START. The data dependences (as they would be exposed in the DAG representation of the program [1]) are analyzed through the processing of VARLIST. These assignments synchronize themselves through the data-flow. The data-driven execution automatically exploits the available instruction level parallelism.

All assignments evaluate the right-hand side (RHS) or source expression. This evaluation results in a pointer to a CDFG object's output (or pseudo-object as defined below). For integer assignments, the left-hand side (LHS) variable or destination is combined with the RHS result object to form a new pair {LHS, result(RHS)} which is added to the front of VARLIST.

The simplest statement is a constant assigned to an integer:⁶

```
a = 5;
```

It doesn't change the CDFG, but adds {a, 5} to the front of VARLIST. The constant 5 is a "pseudo-object" which only holds the value, but does not refer to a CDFG object. Now VARDEF(a) equals 5 at subsequent program points before a is redefined.

Integer assignments can also combine variables already defined and constants:

```
b = a * 2 + 3;
```

In the AST, the RHS is already converted to an expression tree. This tree is transformed to a combination of old and new CDFG objects (which are added to the CDFG) as follows: Each operator (internal node) of the tree is substituted by an ALU with the opcode corresponding to the operator in the tree. If a leaf node is a constant, the ALU's input is directly connected to that constant. If a leaf node is an integer variable var, it is looked up in VARLIST, i. e. VARDEF(var) is retrieved. Then VARDEF(var) (an output of an already existing object in CDFG or a constant) is connected to the ALU's input. The output of the ALU corresponding to the root operator in the expression tree is defined as the *result* of the RHS. Finally,

⁵This method of using a VARLIST is adapted from the Transmogrifier C compiler [3].

⁶Note that we use C syntax for the following examples

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 8

a new pair {LHS, result(RHS)} is added to VARLIST. If the two assignments above are processed, the CDFG with two ALUs in Fig. 2 is created.⁷ Outputs occurring in VARLIST are labeled by Roman numbers. After these two assignments, VARLIST = [{b, I}, {a, 5}]. (The front of the list is on the left side.) Note that all inputs connected to a constant (whether direct from the expression tree or retrieved from VARLIST) must be defined as constant. Inputs defined as constants have a small c next to the input arrow in Fig. 2.

5.2.2 Conditional Integer Assignments

For conditional if-then-else statements containing only integer assignments, objects for condition evaluation are created first. The object event output indicating the condition result is kept for choosing the correct branch result later. Next, both branches are processed in parallel, using separate copies VARLIST1 and VARLIST2 of VARLIST. (VARLIST itself is not changed.) Finally, for all variables added to VARLIST1 or VARLIST2, a new entry for VARLIST is created (combination phase). The valid definitions from VARLIST1 and VARLIST2 are combined with a MUX function, and the correct input is selected by the condition result. For variables only defined in one of the two branches, the multiplexer uses the result retrieved from the original VARLIST for the other branch. If the original VARLIST does not have an entry for this variable, a special "undefined" constant value is used. However, in a functionally correct program this value will never be used. As an optimization, only variables *live* [1] after the if-then-else structure need to be added to VARLIST in the combination phase.

Consider the following example:

```
i = 7;
a = 3;
if (i < 10) {
    a = 5;
    c = 7;
}
else {
    c = a - 1;
    d = 0;
}
```

Fig. 3 shows the resulting CDFG. Before the if-then-else construct, VARLIST = [{a, 3}, {i, 7}]. After processing the branches, for the then branch, VARLIST1 = [{c, 7}, {a, 5}, {a, 3}, {i, 7}], and for the else branch, VARLIST2 = [{d, 0}, {c, I}, {a, 3}, {i, 7}]. After combination, VARLIST = [{d, II}, {c, III}, {a, IV}, {a, 3}, {i, 7}].

Note that case- or switch-statements can be processed, too, since they can – without loss of generality – be converted to nested if-then-else statements.

This processing of conditional statements doesn't need explicit control, either. Both branches are executed in parallel and synchronized by the data-flow.

⁷Note that the input and output names can be deduced from their position, cf. Fig. 1. Also note that the compiler frontend would normally have substituted the second assignment by $b = 13$ (constant propagation). For the simplicity of this explanation, no frontend optimizations are considered in this and the following example.

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 9

5.2.3 Array Accesses

In contrast to the above sections, array accesses have to be controlled explicitly to maintain the correct execution order. For a read access the read address is connected to data input RD. For a write access, the write address is connected to data input WR and the write value to input IN. Alls these inputs are connected to their respective sources through a GATE controlled by START. A STOP event signalling completion of the array access must be assigned to the START variable. Since there's only one START event packet available, only one array access can occur at a time, and the execution order of the original program is maintained. This scheduling scheme is similar to a *one-hot controller* for digital hardware.

If a RAM is read and written at only one program point, the ERD or EWR outputs can be used as STOP events. However, if several read or several write accesses (from different program points) to the same RAM occur, each access produces a ERD or EWR event, respectively. But a STOP event should only be executed for the program point currently executed, the *current access*. This is achieved by connecting the START signals (i.e. those connected to the GATES) of all *other* accesses with the *inverted* START signal of the current access. The resulting signal produces an event for every access, but only for the current access a 1-event. This event is combined (ECOMB) with the RAM's ERD or EWR access. The ECOMB's output will only occur after the access is completed. Because ECOMB OR-combines its event packets, only the current access produces a 1-event. Next, this event is filtered with a 1-FILTER and changed by a 0-CONSTANT, resulting in a STOP signal which produces a 0-event only after the current access is completed as required. See below for an example.

For computing the RAM addresses, the compiler frontend's standard transformation for array accesses can be used. The only difference is that the offset with respect to the RDFP RAM (as determined in the initial mapping phase) must be used.

For several accesses, several sources can be connected to the RD, WR and IN inputs of a RAM. This disables the self synchronization. However, since only one access at a time can happen, the GATES only allow one data packet to arrive at the inputs.

For read accesses, the packets at the OUT output face the same problem as the ERD event packets: They occur for every read access, but must only be used (and forwarded to subsequent operators) for the current access. This can be achieved by connecting the OUT output via a DEMUX function. The Y output of the DEMUX is used, and the X output is left unconnected. The it acts as a selective gate which only forwards packets if its SEL input receives a 1-event, and discards its data input if SEL receives a 0-event. The signal created by the ECOMB described above for the STOP signal creates a 1-event for the current access, and a 0-event otherwise. Using it as the SEL input achieves exactly the desired functionality.

To avoid redundant read accesses, RAM reads are also registered in VARLIST. Instead of an integer variable, an array element is used as first element of the pair. However, a change in a variable occurring in an array index invalidates the information in VARLIST. It must then be removed from it.

The following example shows two read accesses:

```
x = a[i];
y = a[j];
z = x + y;
```

Fig. 4 shows the resulting CDFG. Inputs START (old), i and j should be substituted by the actual functions resulting from the program before the array reads. The signal indicating the STOP of the first access

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 10

is marked by STOP1. Write accesses use the same control events, but instead of one GATE per access for the RD inputs, one GATE for WR and one gate for IN (with the same E input) are used. Also no outputs need to be handled.

Fig. 5 shows the access $a[i] = x$; for the simple case that the RAM is only written once, i. e. at one program point.

This scheme executes RAM accesses correctly, but not very fast since all accesses are synchronized even if this is not necessary. The following optimizations are possible:

- Only accesses to the same RAM are synchronized. Accesses to different arrays can occur concurrently or even in changed order. When there is a data dependency, the accesses self-synchronize automatically. This can be achieved by maintaining a separate START signal for every RAM. At the end of a basic block [1], all these START signals must be combined by a ECOMB to provide a new signal for the next basic block.
- For sequences of either read accesses or write accesses (not mixed) within a basic block, it is possible to stream data into the RAM rather than waiting for the previous access to complete. For this purpose, a combination of MERGE functions selects the RD or WR and IN inputs in the order dictated by the sequence. The MERGEs must be controlled by iterative ESEQs guaranteeing that the inputs are only forwarded in this order. Then only the first access in the sequence needs to be controlled by GATES, the other GATES can be removed to increase throughput. Similarly, the OUT outputs of a read access can be distributed more efficiently for a sequence. A combination of DEMUX functions with the same ESEQ control can be used. For read accesses, the generation of the last output can be sent through a GATE (without the E input connected), thereby producing a STOP event.

Fig. 6 shows the following three array reads in the optimized fashion.

```
x = a[i];
y = a[j];
z = a[k];
```

5.2.4 Input and Output Ports

Input and output ports are processed similar to vector accesses. A read from an input port is like an array read without an address. The input data packet is sent to DEMUX functions which send it to the correct subsequent operators. The STOP signal is generated in the same way as described above for RAM accesses by combining the INPORT's U output with the current and other START signals.

Output ports control the data packets by GATES like array write accesses. The STOP signal is also created as for RAM accesses.

5.2.5 General Conditional Statements

Conditional statements containing either array accesses or inner loops cannot be processed as described in Section 5.2.2. Data packets must only be sent to the active branch. Therefore, a dataflow analysis is

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 11

required. *Used sets* and *defined sets* [1] of both branches must be computed. For all variables in either of these sets DEMUX functions controlled by the IF condition are inserted. They route data packets only to the selected branch. New lists VARLIST1 and VARLIST2 are compiled with the respective outputs of these DEMUX functions. The then-branch is processed with VARLIST1, and the else branch with VARLIST2. Finally, the output values are combined. Since only one branch is ever activated there will not be a conflict due to two packets arriving simultaneously. The combinations will be added to VARLIST after the conditional statement.

5.2.6 FOR Loops

A FOR loop is controlled by a counter CNT. The lower bound (LB), upper bound (UB), and increment (INC) expressions are evaluated like any expressions (see Sections 5.2.1 and 5.2.3) and connected to the respective inputs. The START input is connected to the START signal. The new START signal (after loop execution) is CNT's W output sent through a 1-FILTER and 0-CONSTANT. (W only outputs a 1-event after the counter has terminated.) CNT's V output produces one 0-event for each loop iteration and is therefore used as START for the loop body. Finally, CNT's NEXT input is connected to the START signal at the end of the loop body (i.e. its STOP signal.) This assures that one iteration only starts after the previous one has finished. CNT's X output provides the current value of the loop index variable. For FOR loops, dataflow analysis is required, too.

For all variables *defined* in the loop body and *live* at its beginning, a combination of the input value (from VARLIST at loop entry) and a *feedback value* from the end of the loop is created. Next, each one of these signals is connected to a DEMUX which is controlled by CNT's W output. It sends the input or feedback values back to the loop body (0-event) during loop execution. The VARLIST used in the loop body contains these DEMUX outputs. After loop termination, the input of feedback values are sent to the output of the loop (1-event). The varlist at the end of the loop contains these DEMUX outputs. Inputs not defined in the loop are taken from the input VARLIST.

The processing of the loop body requires some special consideration. Data packets from variables defined outside the loop but only used inside the loop (not redefined) do not lead to the creation of a feedback signal as explained above. Therefore only one packet is available (unless it is a constant), but it is consumed in each loop operation. This would stall the loop operation from the second iteration onwards. Thus it is necessary to multiply the packet for each loop operation. This is achieved by a MDATA function with the SEL input connected to CNT's U output.

These methods allow to process arbitrarily nested loops and conditional statements.

Fig. 7 shows the generated CDFG for the following for loop.

```
a = b + c;
for (i=0; i<=10; i++) {
    a = a + i;
    x[i] = k;
}
```

Note that only one data packet arrives for variables b, c and k, and one final packet is produced for a (out). No GATES are inserted for the RAM write accesses since the packet generation is controlled by the counter anyway.

A Method for Compiling High-Level Language Programs to a Reconfigurable Data-Flow Processor 12

5.2.7 WHILE Loops

WHILE loops are processed similarly. The STOP signal (new START signal) is generated from the loop condition, fed through a 0-FILTER. When the loop finishes, an additional signal (similar to the CNT's W output) must be generated which controls the DEMUXes to generate an output.

5.2.8 Parallelization, Vectorization and Pipelining

The method described so far generates CDFGs performing the HLL program's functionality on an RDFP. However, the program execution is unnecessarily sequentialized by the START signals. In many cases this is too restrictive. Several optimizations are possible.

Independent loops (operating on different variables and arrays) need not be sequentialized. They can use the same START signal, and operate independently. After execution, their STOP signals must be combined by ECOMB, forming a new START signal for the subsequent program parts.

In some cases, loops can be vectorized. This means that loop iterations can overlap, leading to a pipelined data-flow through the operators of the loop body [4]. This technique can be easily applied to the method described here. For FOR loops, the CNT's NEXT input is removed so that CNT counts continuously, thereby overlapping the loop iterations. Since vectorizable loops have no memory access conflicts, the read and write accesses to the same RAM can also overlap. Especially for dual-ported RAM this leads to considerable performance improvements. In this case separate START signals must not only be maintained for each RAM, but also separately for read and write accesses.

Finally, loop transformations like loop unrolling, loop distribution, loop tiling or loop merging [4] can be applied to increase the parallelism and improve performance.

References

- [1] A.V. Aho, R. Sethi, and J.D. Ullman. *Compilers — Principles, Techniques, and Tools*. Addison-Wesley, 1986.
- [2] The Stanford SUIF Compiler Group. Homepage <http://suif.stanford.edu>.
- [3] D. Galloway. The transmogifier C hardware description language and compiler for FPGAs. In *Proc. FPGAs for Custom Computing Machines*, pages 136–144. IEEE Computer Society Press, 1995.
- [4] M. Weinhardt and W. Luk. Pipeline vectorization. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, February 2001.

Akte: PACT27/EP

Applicant: Pact Informationstechnologie GmbH
Leopoldstr. 236
80807 München

Representative:

European Patent Attorney
Claus Peter Pietruk
Heinrich-Lilienfein-Weg 5
D-76229 Karlsruhe

CLAIMS

1. A method for partitioning large computer programs and or algorithms at least part of which is to be executed by an array of reconfigurable units such as ALUS, comprising the steps of
defining a maximum allowable size to be mapped onto the array,
partitioning the program such that its separate parts minimize the overall execution time and providing a mapping onto the array not exceeding the maximum allowable size.
2. A device for partitioning large computer programs and or algorithms at least part of which is to be executed by an array of reconfigurable units such as ALUS, comprising
means for defining a maximum allowable size to be mapped onto the array,
means for partitioning the program such that its separate parts minimize the overall execution time and for providing a mapping onto the array not exceeding the maximum allowable size.

PAGE 1-3
DRAWINGS

Fig. 2:

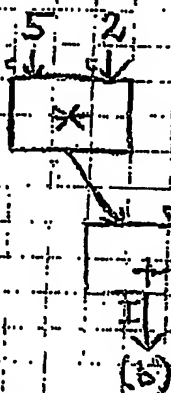


Fig. 3:

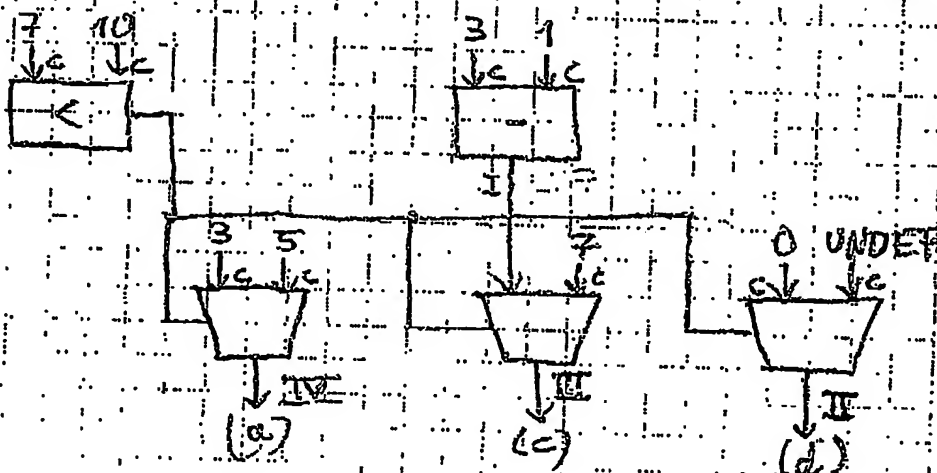
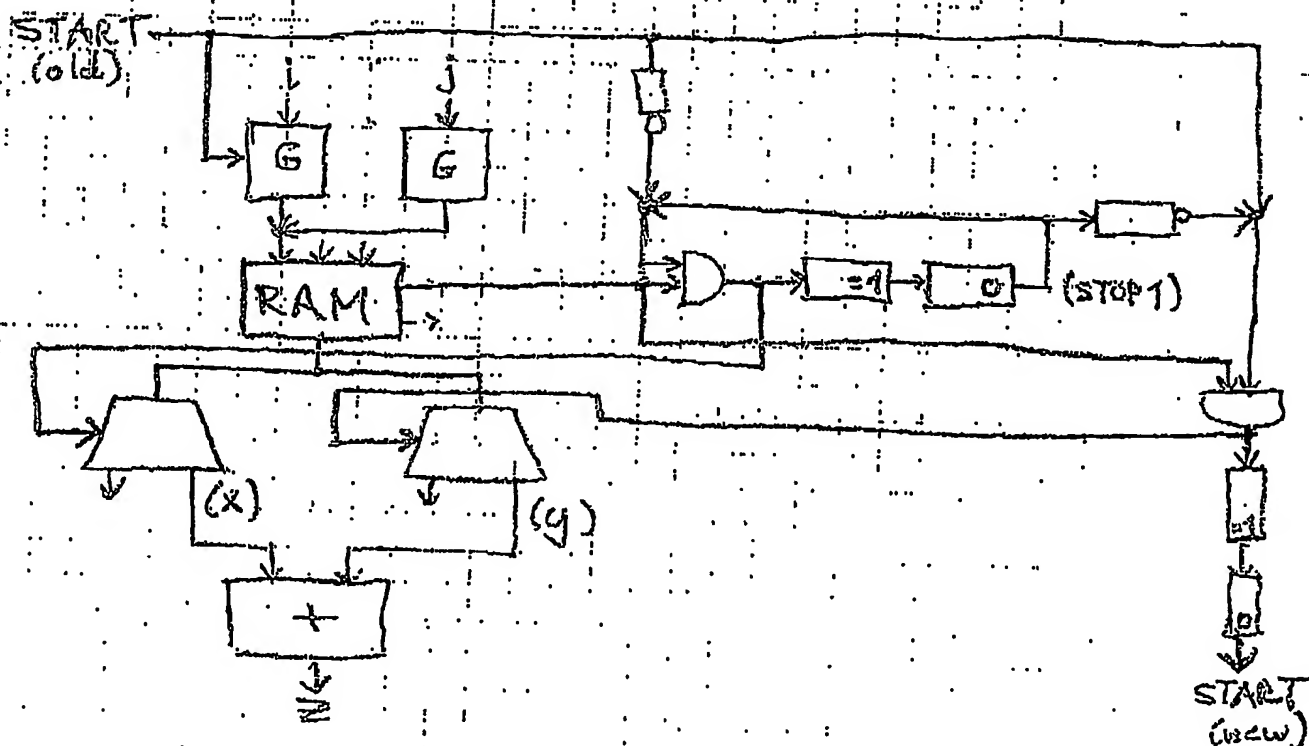


Fig. 4:



PAGE 2-3
DRAWINGS

Fig. 5:

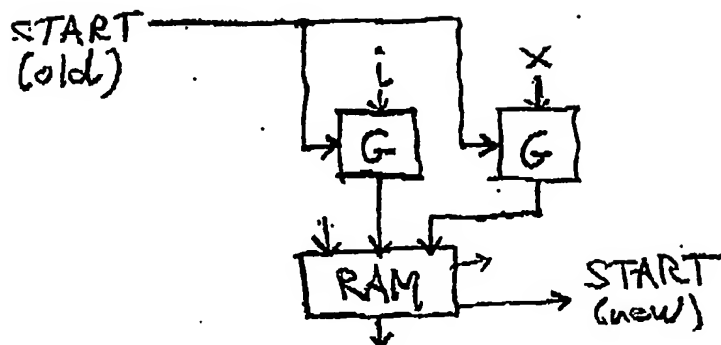


Fig. 6:

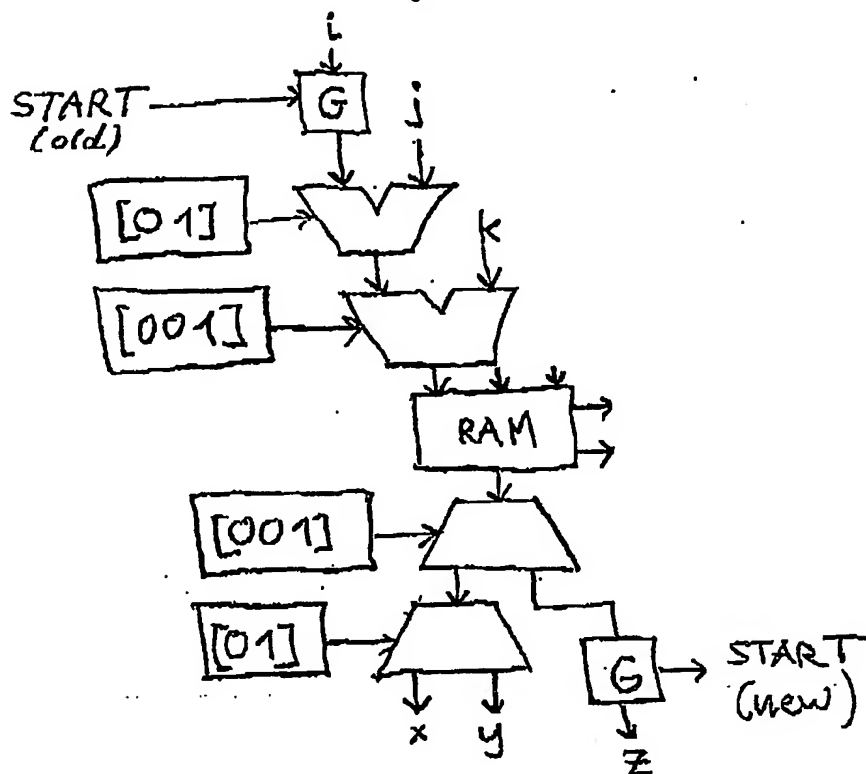


Fig. 7:

START
(old)

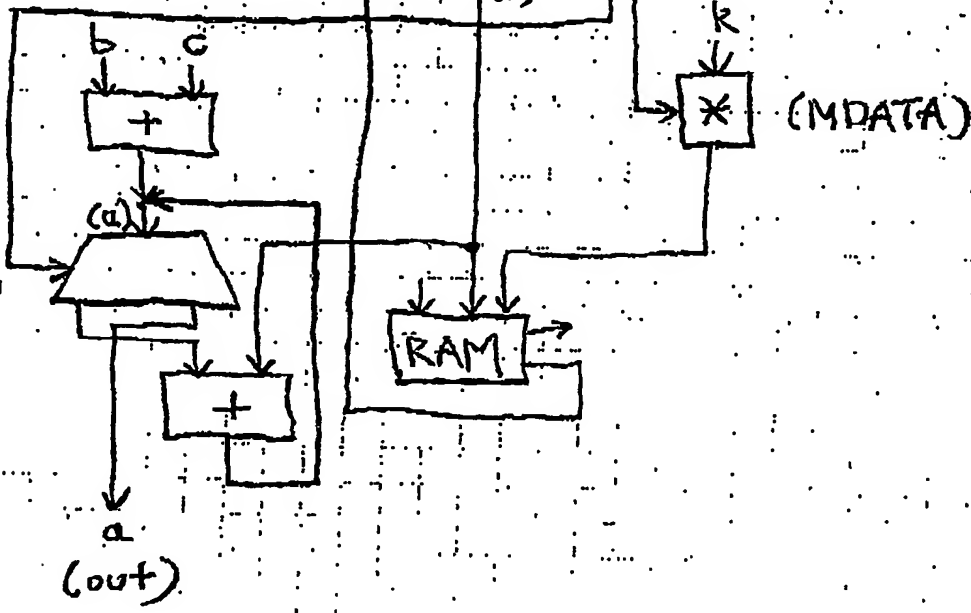
CNT

=1

0

START
(new)

PAGE 3-
DRAWING



**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ BLACK BORDERS
- ☐ IMAGE CUT OFF AT TOP, BOTTOM OR SIDES
- ☐ FADED TEXT OR DRAWING
- ☒ BLURRED OR ILLEGIBLE TEXT OR DRAWING
- ☐ SKEWED/SLANTED IMAGES
- ☐ COLOR OR BLACK AND WHITE PHOTOGRAPHS
- ☐ GRAY SCALE DOCUMENTS
- ☒ LINES OR MARKS ON ORIGINAL DOCUMENT
- ☐ REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY
- ☐ OTHER: _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.